

2018 Machine Learning with R

# Multiple Linear Regression

강 필 성

고려대학교 산업경영공학부

pilsung\_kang@korea.ac.kr

# 목차

I

다중선형회귀분석

II

회귀분석 성능 평가

III

변수 선택

IV

R 실습

# 회귀분석

## ❖ 도요타 코롤라 자동차의 중고차 가격 예측



종속 변수  
(target)

설명 변수  
(attributes, features)

Variable	Description
Price	Offer Price in EUROS
Age_08_04	Age in months as in August 2004
KM	Accumulated Kilometers on odometer
Fuel_Type	Fuel Type (Petrol, Diesel, CNG)
HP	Horse Power
Met_Color	Metallic Color? (Yes=1, No=0)
Automatic	Automatic (Yes=1, No=0)
CC	Cylinder Volume in cubic centimeters
Doors	Number of doors
Quarterly_Tax	Quarterly road tax in EUROS
Weight	Weight in Kilograms

# 다중회귀분석: 목적

## ❖ 목적

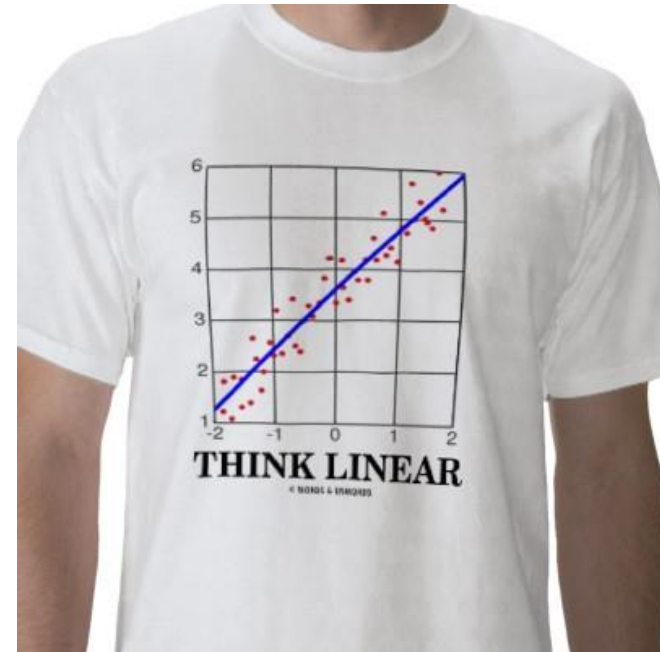
- 종속변수  $Y$ 와 설명변수 집합  $X_1, X_2, \dots, X_p$ 사이의 관계를 **선형으로 가정**하고 이를 가장 잘 설명할 수 있는 회귀 계수(regression coefficients)를 추정

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \cdots + \beta_d x_d + \epsilon$$

unexplained

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \cdots + \hat{\beta}_d x_d$$

coefficients



# 다중회귀분석: 탐색적 vs. 예측적

## ❖ 탐색적(Explanatory) 회귀분석 vs. 예측적(Predictive) 회귀분석

### Explanatory Regression

- Explain relationship between predictors (explanatory variables) and target.
- Familiar use of regression in data analysis.
- Model Goal: Fit the data well and understand the contribution of explanatory variables to the model.
- “goodness-of-fit”:  $R^2$ , residual analysis, p-values.

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$$

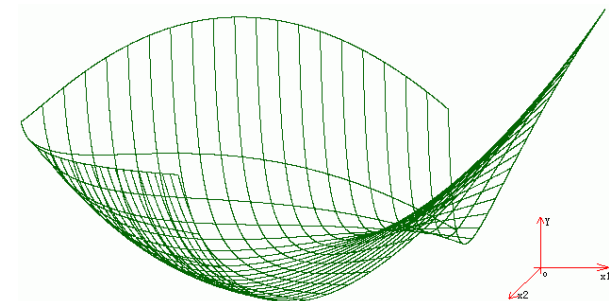
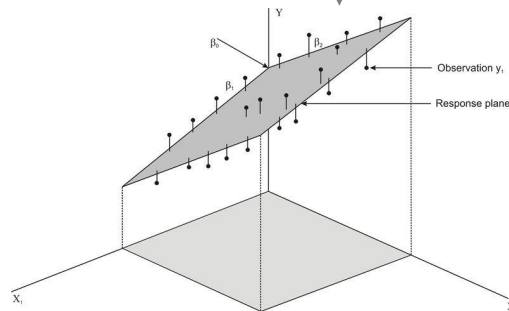
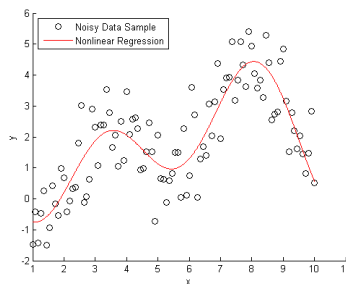
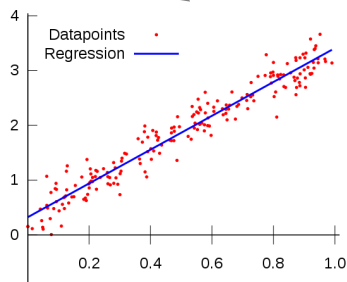
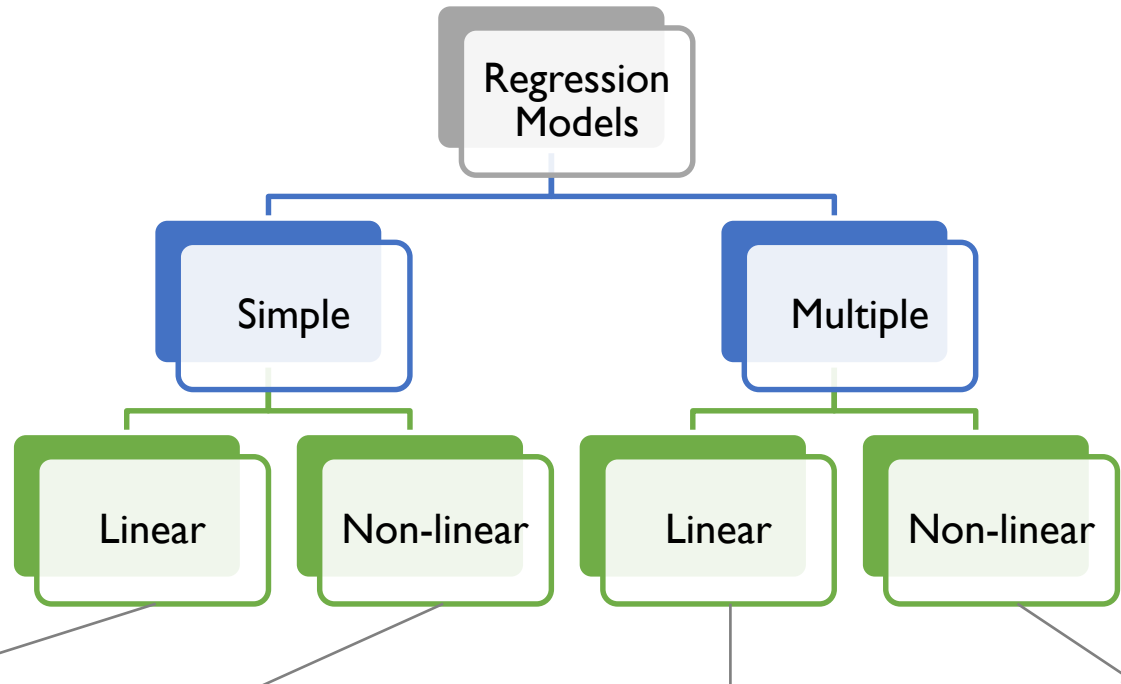
### Predictive Regression

- Predict target values in other data where we have predictor values, but not target values.
- Classic data mining context
- Model Goal: Optimize predictive accuracy
- Train model on training data
- Assess performance on validation (hold-out) data
- Explaining role of predictors is not primary purpose (but useful)

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$$

# 회귀분석의 형태

❖ 변수의 수와 추정되는 함수의 형태에 따라 아래와 같이 구분 가능

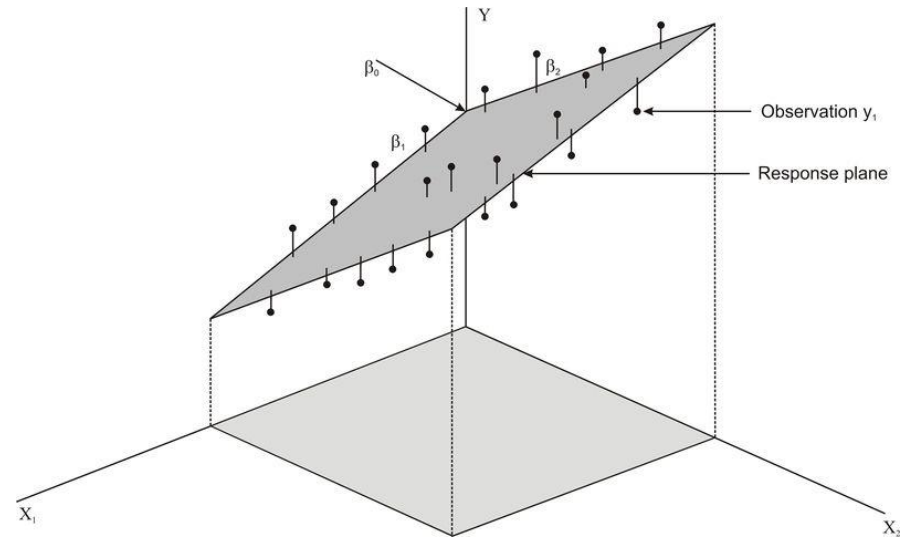
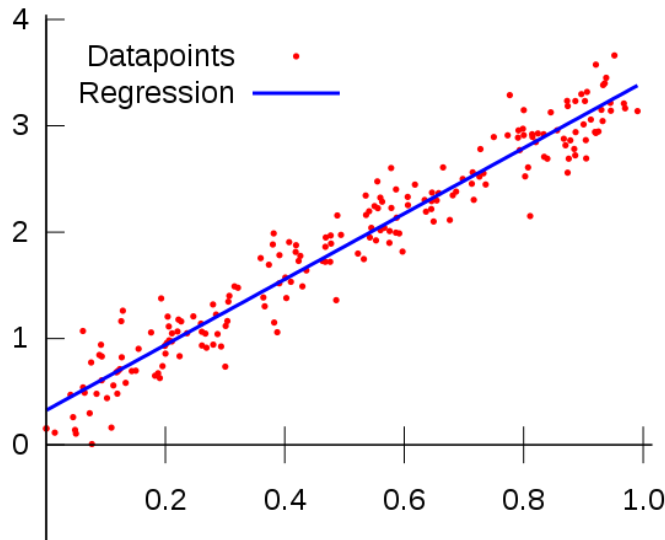


# 선형 회귀분석

## ❖ 선형 회귀 분석: Linear Regression

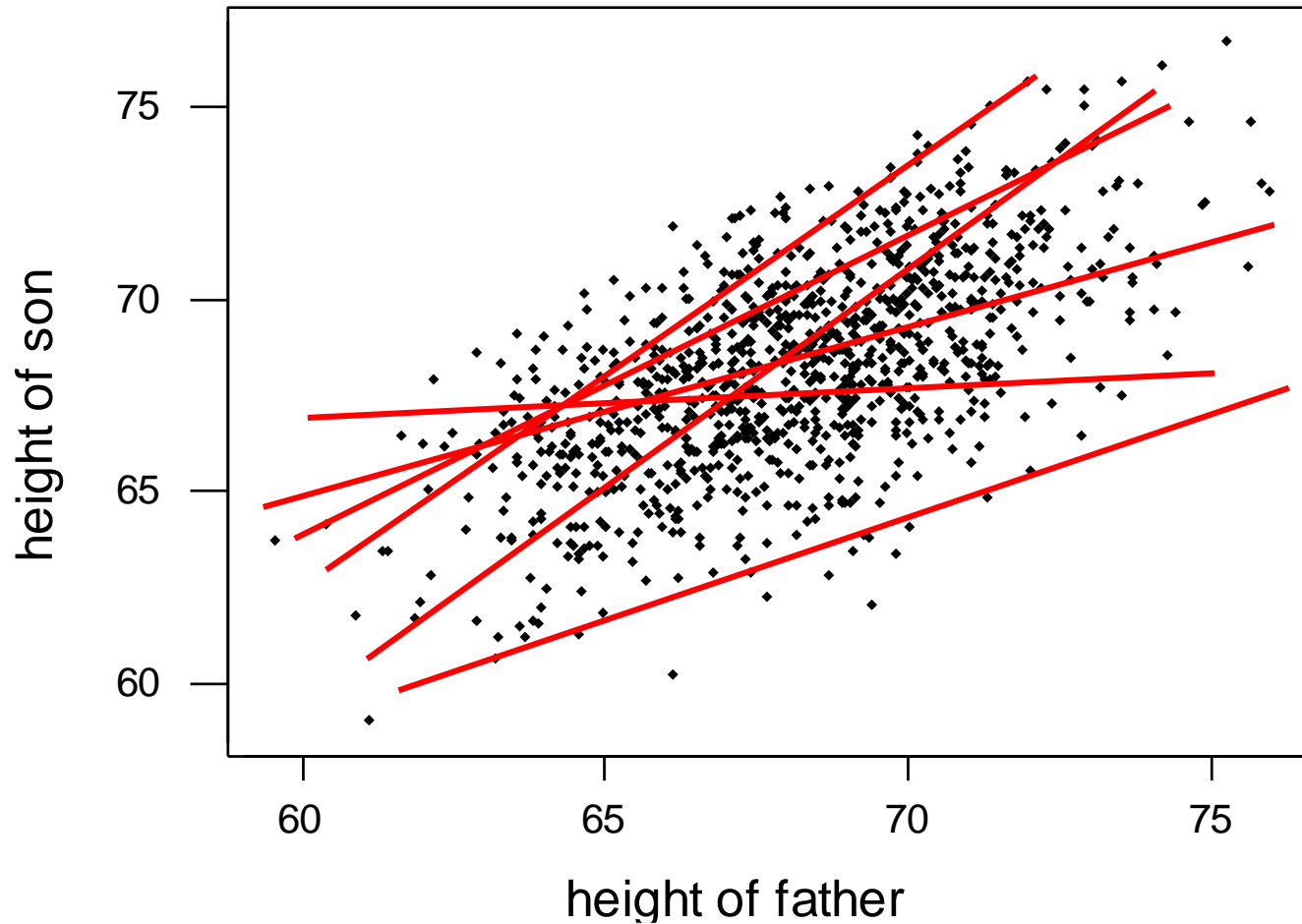
- 반응변수와 설명변수 사이의 관계를 선형으로 표현

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \cdots + \hat{\beta}_d x_d$$



# 선형 회귀분석

❖ 어떤 직선이 설명변수와 종속변수를 가장 잘 표현하는가?

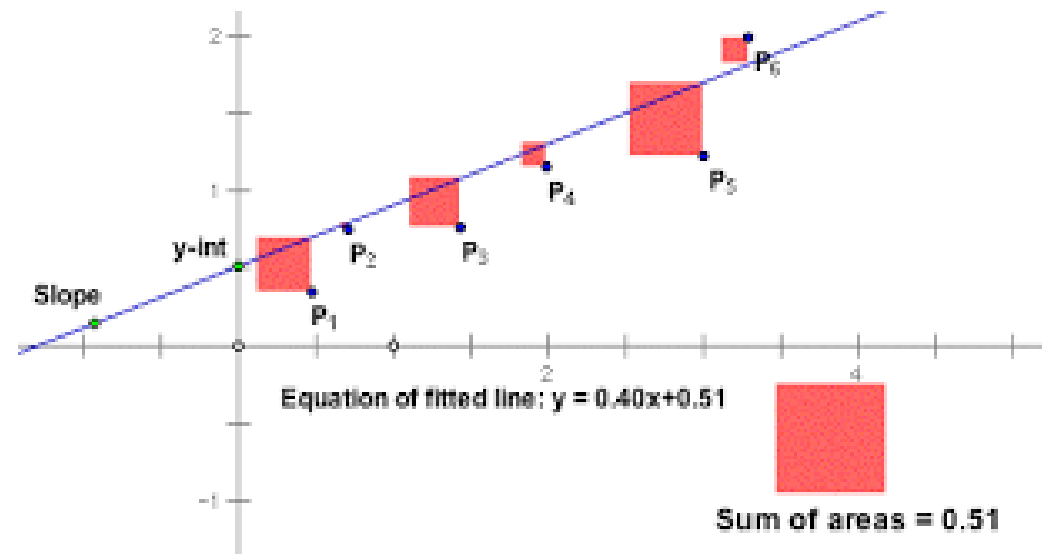
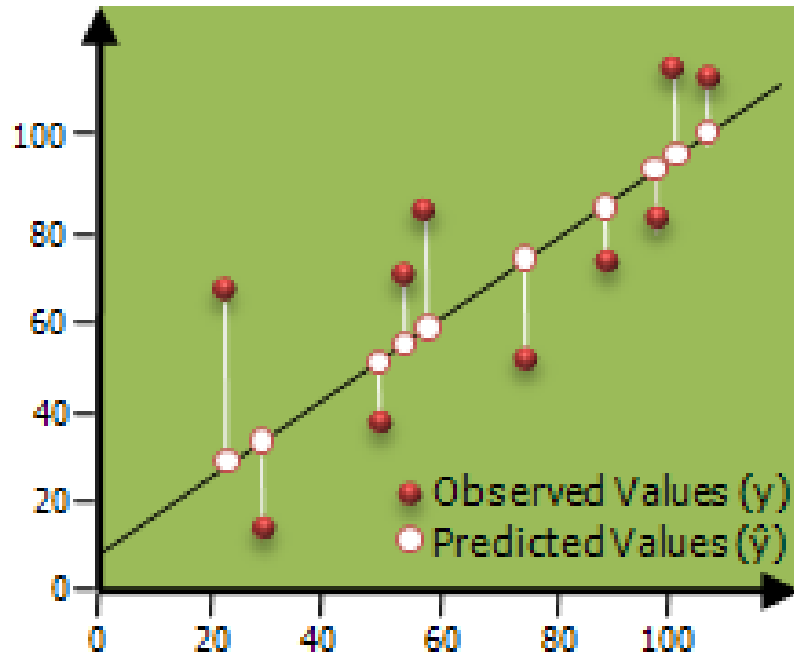




# 선형 회귀분석

## ❖ 최소자승법: Ordinary Least Squares (OLS)

- 추정된 회귀식에 의해 결정된 값과 실제 종속변수 값의 차이를 최소한으로 줄이는 것을 목적으로 함



# 다중회귀분석: 회귀 계수의 추정

## ❖ 회귀 계수의 추정

### ■ 최소자승법: Ordinary least square (OLS)

✓ Actual target:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \cdots + \beta_d x_d + \epsilon$

✓ Predicted target:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \cdots + \hat{\beta}_d x_d$

✓ **목적**: 실제 종속변수 값과 예측된 종속변수 값 사이의 오차 제곱합을 최소화

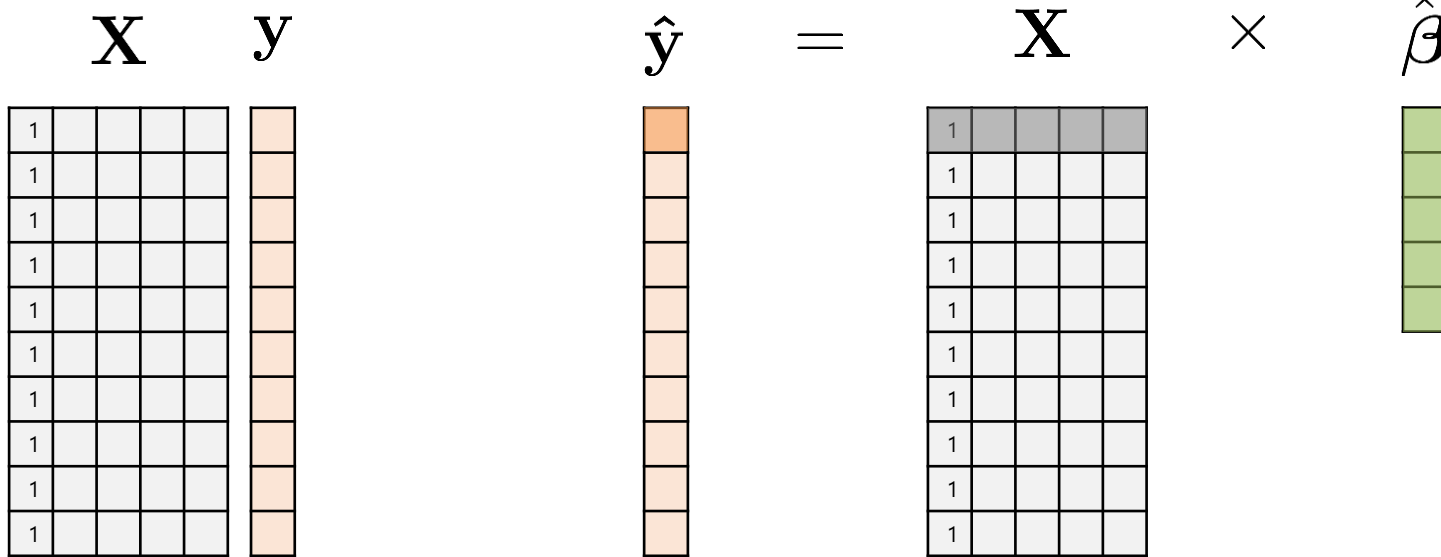
$$\begin{aligned} \min \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ = \frac{1}{2} (y_i - \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} \cdots + \hat{\beta}_d x_{id})^2 \end{aligned}$$

# 다중회귀분석: 회귀 계수의 추정

❖ 최소 자승법: 행렬을 이용한 해 구하기

$\mathbf{X} : n \times (d + 1)$  matrix,  $\mathbf{y} : n \times 1$  vector

$\hat{\boldsymbol{\beta}} : (d + 1) \times 1$  vector



상수항을 취급하기 위한 장치

# 다중회귀분석: 회귀 계수의 추정

❖ 최소 자승법: 행렬을 이용한 해 구하기

$\mathbf{X} : n \times (d + 1)$  matrix,  $\mathbf{y} : n \times 1$  vector

$\hat{\boldsymbol{\beta}} : (d + 1) \times 1$  vector

$$\min E(\mathbf{X}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

$$\Rightarrow \frac{\partial E(\mathbf{X})}{\partial \hat{\boldsymbol{\beta}}} = -\mathbf{X}^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0$$

$$\Rightarrow \mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} = 0$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \longrightarrow \text{학습 데이터에 대해 유일하고 명시적인 해(solution)가 존재!}$$

# 다중회귀분석: 회귀 계수의 추정

❖ 최소 자승법: 행렬을 이용한 해 구하기

$$\hat{\beta} = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\beta} = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

회귀 계수에 대한 Closed form solution이 존재

# 다중회귀분석: 회귀 계수의 추정

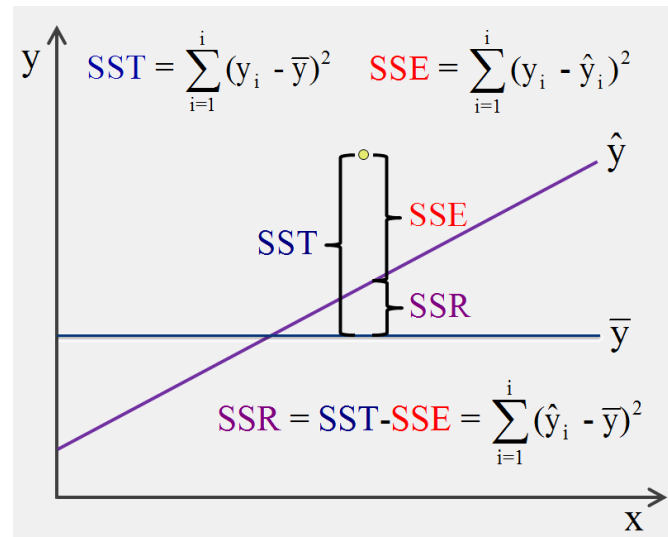
## ❖ 최소자승법

- 아래 조건을 만족할 경우 최소자승법으로 구한 회귀계수  $\beta$ 는 최적해임
  - 오차항  $\varepsilon$  이 정규분포를 따름
  - 설명변수와 종속변수 사이에 선형관계가 성립함
  - 각 관측치들은 서로 독립
  - 종속변수  $Y$ 에 대한 오차항(residual)은 설명변수 값의 범위에 관계없이 일정함(homoskedasticity)

# 다중회귀분석: 회귀 계수의 추정

## ❖ 회귀모형의 적합도

### ■ Sum-of-Squares Decomposition



$$\underbrace{\sum_{j=1}^n (y_j - \bar{y})^2}_{\substack{\text{(total sum of squares)} \\ \text{about mean}}} = \underbrace{\sum_{j=1}^n (\hat{y}_j - \bar{y})^2}_{\substack{\text{(regression)} \\ \text{sum of squares}}} + \underbrace{\sum_{j=1}^n \hat{\varepsilon}_j^2}_{\substack{\text{(residual (error))} \\ \text{sum of squares}}} .$$

**SST**
**SSR**
**SSE**

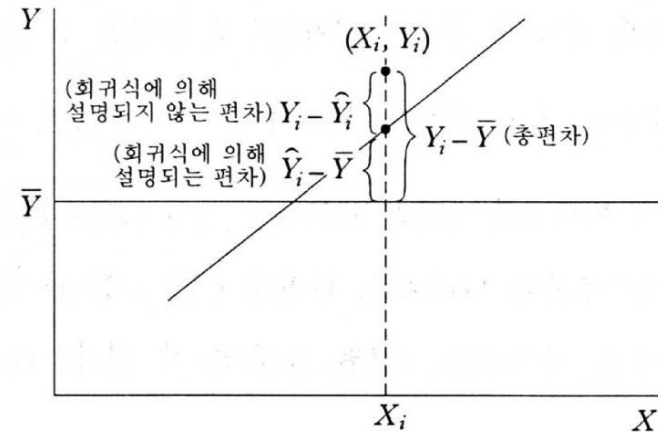
# 다중회귀분석: 회귀 계수의 추정

## ❖ 회귀모형의 적합도

### ■ 결정계수( $R^2$ ):

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SSR}{SST}$$

$$0 \leq R^2 \leq 1$$



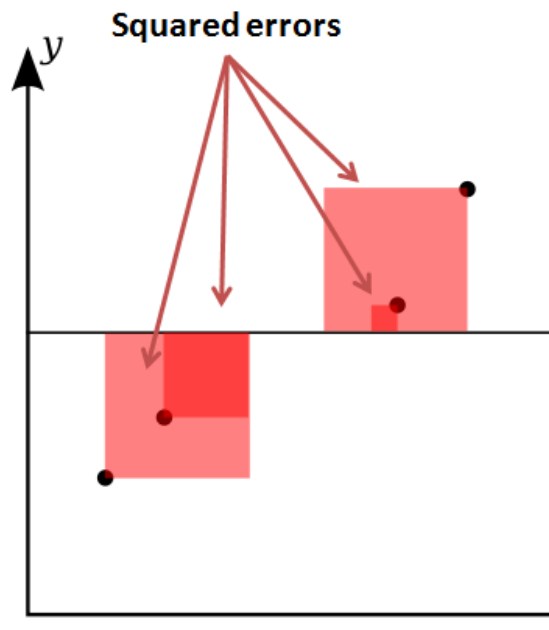
- ✓ 반응변수 ( $Y$ )의 전체 변동 중 예측변수( $X$ )가 차지하는 변동의 비율
- ✓  $R^2$ 는 0과 1 사이에 존재
- ✓  $R^2=1$ : 회귀직선으로  $Y$ 의 총변동이 완전히 설명됨 (모든 측정값들이 회귀직선 위에 있는 경우)
- ✓  $R^2=0$ : 추정된 회귀직선은  $X$ 와  $Y$ 의 관계를 전혀 설명하지 못함



# 다중회귀분석: 회귀 계수의 추정

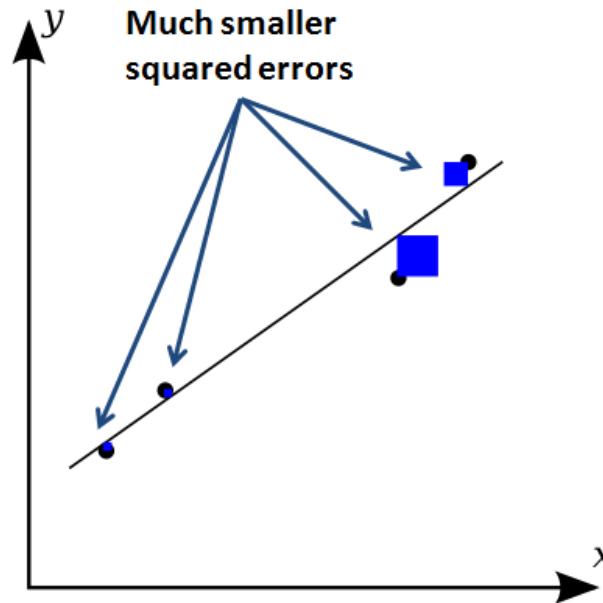
## ❖ 회귀모형의 적합도

### ■ Graphical Interpretation



Computationally:

$$R\text{-squared} = 1 - \left[ \frac{SS_{\text{error}}}{SS_{\text{total}}} \right]$$



Conceptually:

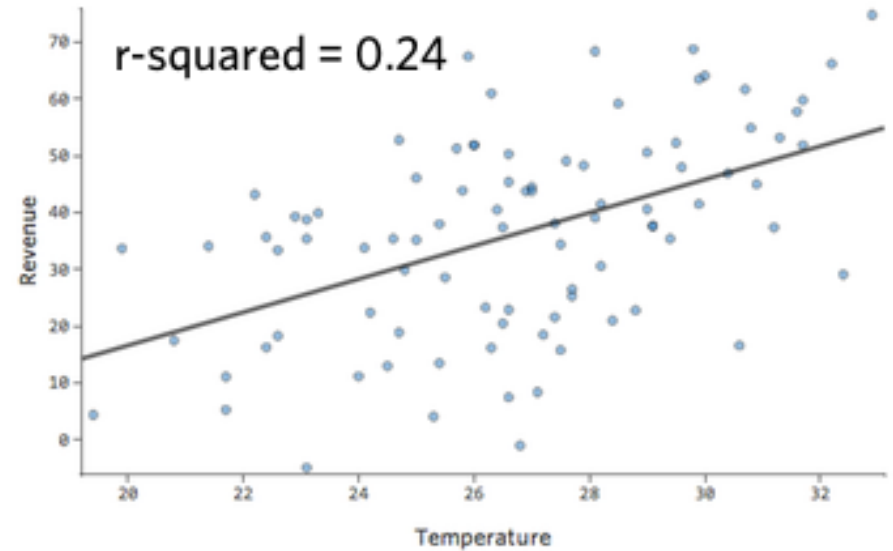
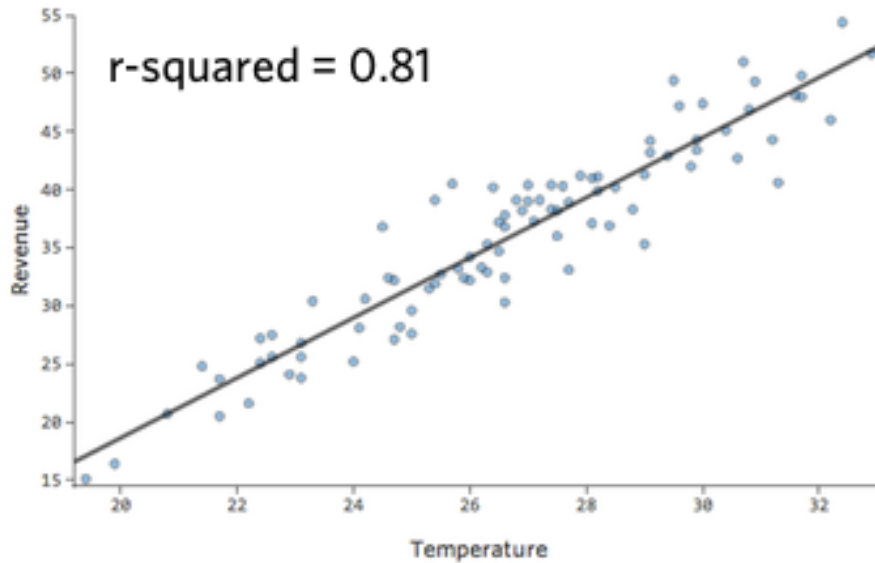
Force x and y to be independent, calculate **the squared error**.

Allow for a relationship between x and y, does this reduce your **error**?

# 다중회귀분석: 회귀 계수의 추정

## ❖ 회귀모형의 적합도

### ■ 결정계수( $R^2$ ):



# 다중회귀분석: 회귀 계수의 추정

## ❖ 회귀모형의 적합도

- 수정 결정계수(Adjusted  $R^2$ ):

$$R_{adj}^2 = 1 - \left[ \frac{n-1}{n-(p+1)} \right] \frac{SSE}{SST} \leq 1 - \frac{SSE}{SST} = R^2$$

- ✓  $R^2$ 는 유의하지 않은 변수가 추가되어도 항상 증가
- ✓ 수정  $R^2$ 는 이러한 단점을 앞에 계수를 곱해줌으로써 보정
- ✓ 유의하지 않은 변수가 추가될 경우 수정 결정계수는 증가하지 않음

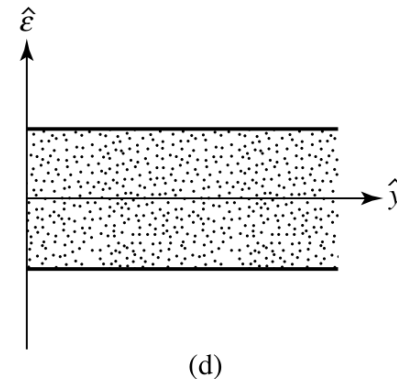
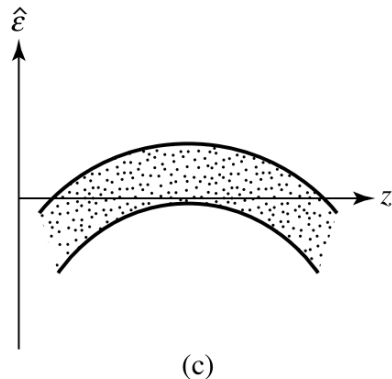
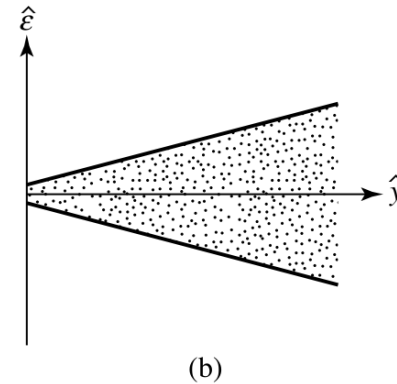
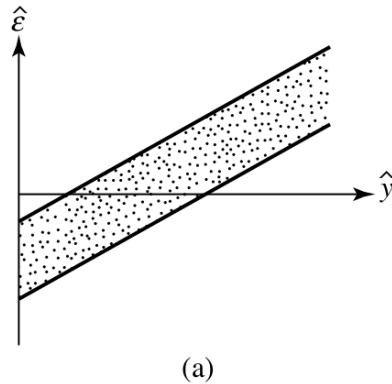
## ❖ 모형의 검토

- 추정된 모형이 다음 가정을 만족하는지 확인
  - ✓ 예측변수와 반응변수 간 관계가 선형
  - ✓ 오차항들이 서로 독립
  - ✓ 오차항은 평균이 0이며 분산이 일정한 정규분포를 따름

# 다중회귀분석: 모형의 적합도 평가

## ❖ Residual Plot:

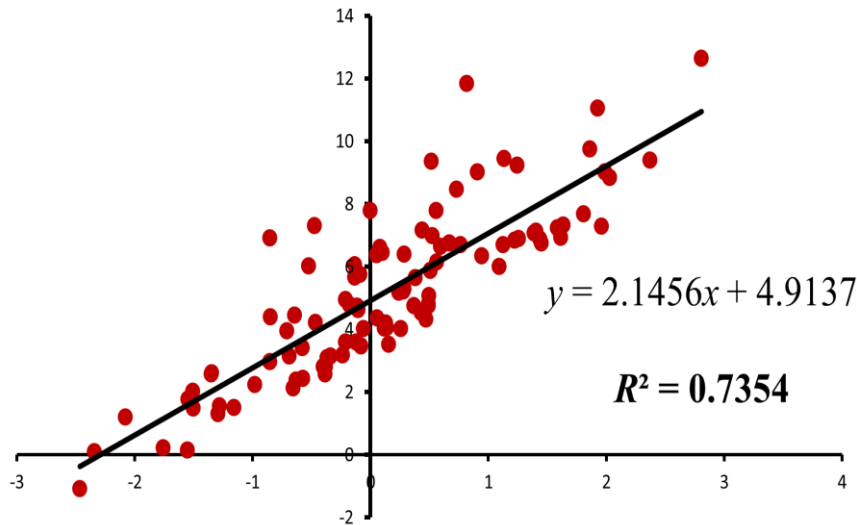
- 종속변수  $Y$ 에 대한 오차항(residual)은 설명변수 값의 범위에 관계없이 일정함(homoskedasticity)을 유지하는지 평가



# 다중회귀분석: 모형의 적합도 평가

## ❖ 잔차의 정규성

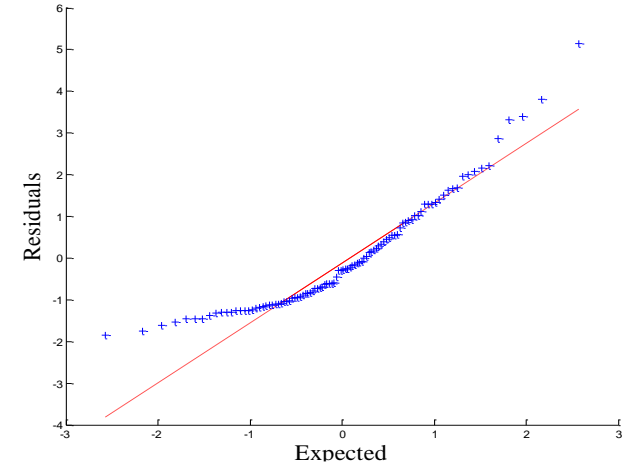
$$y = 2x + \varepsilon, \quad \varepsilon \sim \text{Gamma}(2,1)$$



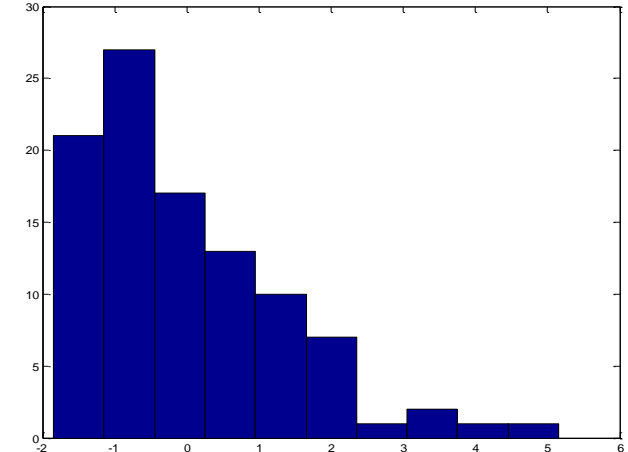
Regression model



QQ Plot of Residuals



Histogram of Residuals



# 다중회귀분석: 예시

## ❖ 예시: 도요타 코롤라 중고차 가격 예측

Y

X

Price	Age_08_04	KM	Fuel_Type	HP	Met_Color	Automatic	cc	Doors	Quarterly_Tax	Weight
13500	23	46986	Diesel	90	1	0	2000	3	210	1165
13750	23	72937	Diesel	90	1	0	2000	3	210	1165
13950	24	41711	Diesel	90	1	0	2000	3	210	1165
14950	26	48000	Diesel	90	0	0	2000	3	210	1165
13750	30	38500	Diesel	90	0	0	2000	3	210	1170
12950	32	61000	Diesel	90	0	0	2000	3	210	1170
16900	27	94612	Diesel	90	1	0	2000	3	210	1245
18600	30	75889	Diesel	90	1	0	2000	3	210	1245
21500	27	19700	Petrol	192	0	0	1800	3	100	1185
12950	23	71138	Diesel	69	0	0	1900	3	185	1105
20950	25	31461	Petrol	192	0	0	1800	3	100	1185
19950	22	43610	Petrol	192	0	0	1800	3	100	1185
19600	25	32189	Petrol	192	0	0	1800	3	100	1185
21500	31	23000	Petrol	192	1	0	1800	3	100	1185
22500	32	34131	Petrol	192	1	0	1800	3	100	1185
22000	28	18739	Petrol	192	0	0	1800	3	100	1185
22750	30	34000	Petrol	192	1	0	1800	3	100	1185
17950	24	21716	Petrol	110	1	0	1600	3	85	1105
16750	24	25563	Petrol	110	0	0	1600	3	19	1065

# 다중회귀분석: 예시

## ❖ 데이터 전처리

- Fuel type 변수에 대한 1-of-C coding 변환

	Fuel_type = Diesel	Fuel_type = Petrol	Fuel_type = CNG
Diesel	1	0	0
Petrol	0	1	0
CNG	0	0	1

## ❖ 데이터 구분

- 학습용 데이터 60%, 검증용 데이터 40%

Id	Model	Price	Age_08_04	Mfg_Month	Mfg_Year	KM	Fuel_Type_Diesel	Fuel_Type_Petrol
1	RRA 2/3-Doors	13500	23	10	2002	46986	1	0
4	RRA 2/3-Doors	14950	26	7	2002	48000	1	0
5	SOL 2/3-Doors	13750	30	3	2002	38500	1	0
6	SOL 2/3-Doors	12950	32	1	2002	61000	1	0
9	VT I 2/3-Doors	21500	27	6	2002	19700	0	1
10	RRA 2/3-Doors	12950	23	10	2002	71138	1	0
12	BNS 2/3-Doors	19950	22	11	2002	43610	0	1
17	ORT 2/3-Doors	22750	30	3	2002	34000	0	1

# 다중회귀분석: 예시

## ❖ 다중회귀분석 결과물 해석

- 다중회귀분석을 수행하고 나면 다음과 같은 표를 결과로 얻을 수 있음

Input variables	Coefficient	Std. Error	p-value	SS
Constant term	-3608.418457	1458.620728	0.0137	97276410000
Age_08_04	-123.8319168	3.367589	0	8033339000
KM	-0.017482	0.00175105	0	251574500
Fuel_Type_Diesel	210.9862518	474.9978333	0.6571036	6212673
Fuel_Type_Petrol	2522.066895	463.6594238	0.00000008	4594.9375
HP	20.71352959	4.67398977	0.00001152	330138600
Met_Color	-50.48505402	97.85591125	0.60614568	596053.75
Automatic	178.1519013	212.0528565	0.40124047	19223190
cc	0.01385481	0.09319961	0.88188446	1272449
Doors	20.02487946	51.0899086	0.69526076	39265060
Quarterly_Tax	16.7742424	2.09381151	0	160667200
Weight	15.41666317	1.40446579	0	214696000



# 다중회귀분석: 예시

## ❖ 다중회귀분석 결과물 해석

### ■ 회귀계수: Coefficient

- ✓ 선형회귀분석에서 각 변수에 대응하는 베타값임
- ✓ 해당 변수가 1단위 증가할 때 종속변수의 변화량을 의미
- ✓ 양수이면 해당 설명변수와 종속변수는 양의 상관관계, 음수이면 음의 상관관계

Input variables	Coefficient	Std. Error	p-value	SS
Constant term	-3608.418457	1458.620728	0.0137	97276410000
Age_08_04	-123.8319168	3.367589	0	8033339000
KM	-0.017482	0.00175105	0	251574500
Fuel_Type_Diesel	210.9862518	474.9978333	0.6571036	6212673
Fuel_Type_Petrol	2522.066895	463.6594238	0.00000008	4594.9375
HP	20.71352959	4.67398977	0.00001152	330138600
Met_Color	-50.48505402	97.85591125	0.60614568	596053.75
Automatic	178.1519013	212.0528565	0.40124047	19223190
cc	0.01385481	0.09319961	0.88188446	1272449
Doors	20.02487946	51.0899086	0.69526076	39265060
Quarterly_Tax	16.7742424	2.09381151	0	160667200
Weight	15.41666317	1.40446579	0	214696000

# 다중회귀분석: 예시

## ❖ 다중회귀분석 결과물 해석

### ■ 유의확률: p-value

- ✓ 선형회귀분석에서 해당 변수가 통계적으로 유의미한지 알려주는 지표
- ✓ 0에 가까울수록 모델링에 중요한 변수이며, 1에 가까울수록 유의미하지 않은 변수임
- ✓ 특정 유의수준( $\alpha$ )을 설정하여 해당 값 미만의 변수만을 사용하여 다시 선형회귀분석을 구축하는 것도 가능함 (주로  $\alpha = 0.05$  사용)

Input variables	Coefficient	Std. Error	p-value	SS
Constant term	-3608.418457	1458.620728	0.0137	97276410000
Age_08_04	-123.8319168	3.367589	0	8033339000
KM	-0.017482	0.00175105	0	251574500
Fuel_Type_Diesel	210.9862518	474.9978333	0.6571036	6212673
Fuel_Type_Petrol	2522.066895	463.6594238	0.00000008	4594.9375
HP	20.71352959	4.67398977	0.00001152	330138600
Met_Color	-50.48505402	97.85591125	0.60614568	596053.75
Automatic	178.1519013	212.0528565	0.40124047	19223190
cc	0.01385481	0.09319961	0.88188446	1272449
Doors	20.02487946	51.0899086	0.69526076	39265060
Quarterly_Tax	16.7742424	2.09381151	0	160667200
Weight	15.41666317	1.40446579	0	214696000

# 다중회귀분석: 예시

## ❖ 실제 종속변수 값과 예측된 종속변수 값의 차이(잔차) 분석

Predicted Value	Actual Value	Residual
15863.86944	13750	-2113.869439
16285.93045	13950	-2335.930454
16222.95248	16900	677.047525
16178.77221	18600	2421.227789
19276.03039	20950	1673.969611
19263.30349	19600	336.6965066
18630.46904	21500	2869.530964
18312.04498	22500	4187.955022
19126.94064	22000	2873.059357
16808.77828	16950	141.2217206
15885.80362	16950	1064.196384
15873.97887	16250	376.0211263
15601.22471	15750	148.7752903
15476.63164	15950	473.3683568
15544.83584	14950	-594.835836
15562.25552	14750	-812.2555172
15222.12869	16750	1527.871313
17782.33234	19000	1217.667664

# 목차

I

다중선형회귀분석

II

회귀분석 성능 평가

III

변수 선택

IV

R 실습

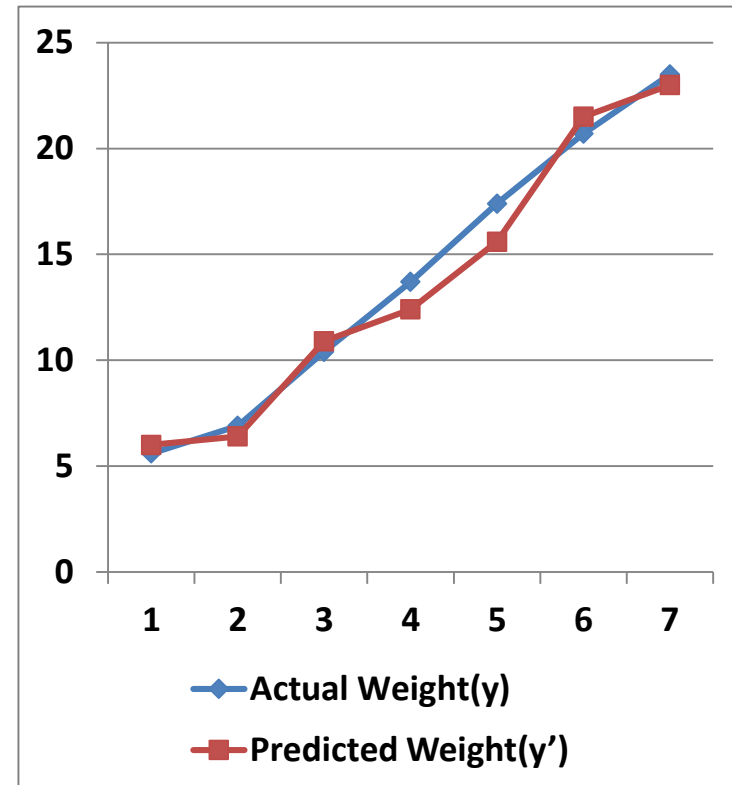
# 회귀분석 성능평가

I

## 예시

- 나이에 따른 아기의 몸무게 예측

Age	Actual Weight(y)	Predicted Weight(y')
1	5.6	6.0
2	6.9	6.4
3	10.4	10.9
4	13.7	12.4
5	17.4	15.6
6	20.7	21.5
7	23.5	23.0



# 회귀분석 성능평가

2

## 성능지표 1: 평균오차

- 실제 값에 비해 과대/과소 추정 여부를 판단
- 부호로 인해 잘못된 결론을 내릴 위험이 있음

$$\begin{aligned} \text{Average error} &= \frac{1}{n} \sum_{i=1}^n (y - y') \\ &= 0.342 \end{aligned}$$

Age	Actual Weight(y)	Predicted Weight(y')
1	5.6	6.0
2	6.9	6.4
3	10.4	10.9
4	13.7	12.4
5	17.4	15.6
6	20.7	21.5
7	23.5	23.0

# 회귀분석 성능평가

## 성능지표 2: 평균 절대 오차(Mean absolute error; MAE)

- 실제 값과 예측 값 사이의 절대적인 오차의 평균을 이용

$$MAE = \frac{1}{n} \sum_{i=1}^n |y - y'|$$

$$= 0.829$$

Age	Actual Weight(y)	Predicted Weight(y')
1	5.6	6.0
2	6.9	6.4
3	10.4	10.9
4	13.7	12.4
5	17.4	15.6
6	20.7	21.5
7	23.5	23.0

# 회귀분석 성능평가

## 성능지표 3: Mean absolute percentage error (MAPE)

- 실제값 대비 얼마나 예측 값이 차이가 있는지를 %로 표현
- 상대적인 오차를 추정하는데 주로 사용

$$MAPE = 100\% \times \frac{1}{n} \sum_{i=1}^n \frac{|y - y'|}{|y|}$$

$$= 6.43\%$$

Age	Actual Weight(y)	Predicted Weight(y')
1	5.6	6.0
2	6.9	6.4
3	10.4	10.9
4	13.7	12.4
5	17.4	15.6
6	20.7	21.5
7	23.5	23.0



# 회귀분석 성능평가

## 성능지표 4 & 5: (Root) Mean squared error ((R)MSE)

- 부호의 영향을 제거하기 위해 절대값이 아닌 제곱을 취한 지표

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - y')^2$$

$$= 0.926$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y - y')^2}$$

$$= 0.962$$

Age	Actual Weight(y)	Predicted Weight(y')
1	5.6	6.0
2	6.9	6.4
3	10.4	10.9
4	13.7	12.4
5	17.4	15.6
6	20.7	21.5
7	23.5	23.0

# 회귀분석 성능평가

## ❖ 학습 및 검증 데이터에 대한 성능 평가

### Training Data scoring - Summary Report

Total sum of squared errors	RMS Error	Average Error
1514553377	1325.527246	-0.000426154

### Validation Data scoring - Summary Report

Total sum of squared errors	RMS Error	Average Error
1021587500	1334.079894	116.3728779

# 목차

I

다중선형회귀분석

II

회귀분석 성능 평가

III

변수 선택

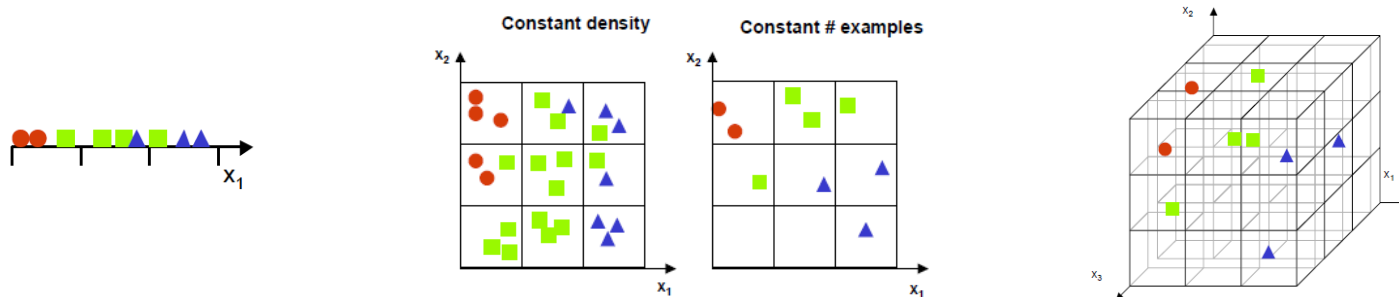
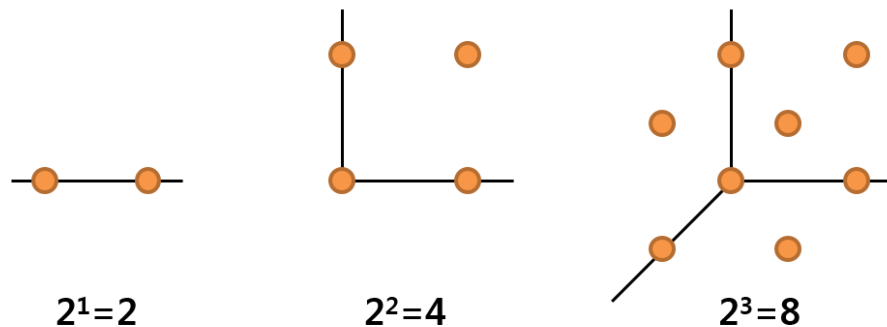
IV

R 실습

# 차원축소: Dimensionality Reduction

## ❖ 차원의 저주 (Curse of Dimensionality)

- 동등한 설명력을 갖기 위해서는 변수가 증가할 때 필요한 개체의 수는 기하급수적으로 증가함



*“If there are various logical ways to explain a certain phenomenon, the simplest is the best” - Occam’s Razor*

# 차원축소: Dimensionality Reduction

## ❖ 차원축소: 배경

- 이론적으로는 변수의 수가 증가할수록 모델의 성능이 향상됨 (변수간 독립성 만족 시)
- 실제 상황에서는 변수간 독립성 가정 위배, 노이즈 존재 등으로 인해 변수의 수가 일정 수준 이상 증가하면 모델의 성능이 저하되는 경향이 있음

## ❖ 차원축소: 목적

- 향후 분석 과정에서 성능을 저하시키지 않는 최소한의 변수 집합 판별

## ❖ 차원축소: 효과

- 변수간 상관성을 제거하여 결과의 통계적 유의성 제고
- 사후 처리(post-processing)의 단순화
- 주요 정보를 보존한 상태에서 중복되거나 불필요한 정보만 제거
- 고차원의 정보를 저차원으로 축소하여 시각화(visualization) 가능

# 차원축소: Dimensionality Reduction

## ❖ 차원축소 방식

### ■ 교사적 차원축소 (Supervised dimensionality reduction)

- ✓ 축소된 차원의 적합성을 검증하는데 있어 데이터마이닝 모델을 적용
- ✓ 동일한 데이터라도 적용되는 데이터마이닝 모델에 따라 축소된 차원의 결과가 달라질 수 있음

### ■ 비교사적 차원축소 (Unsupervised dimensionality reduction)

- ✓ 축소된 차원의 적합성을 검증하는데 있어 데이터마이닝 모델을 적용하지 않음
- ✓ 특정 기법에 따른 차원축소 결과는 동일함

# 차원축소: Dimensionality Reduction

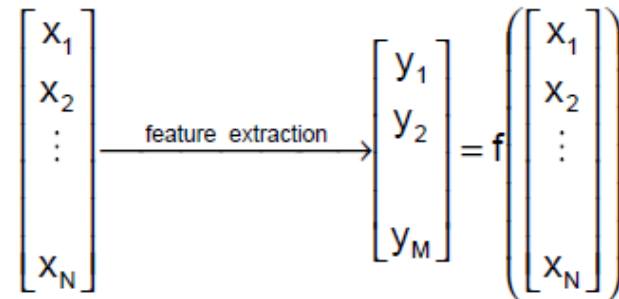
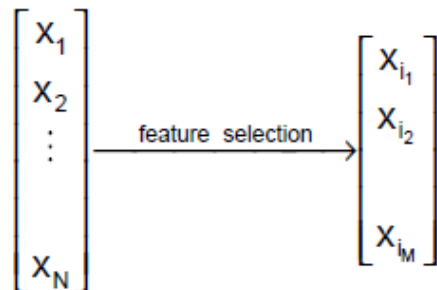
## ❖ 차원축소 기법

### ■ 변수 선택(variable/feature selection)

- ✓ 원래의 변수 집단으로부터 유용할 것으로 판단되는 소수의 변수들을 선택
- ✓ Filter – 변수 선택 과정과 모델 구축 과정이 독립적
- ✓ Wrapper – 변수 선택 과정이 데이터마이닝 모델의 결과를 최적화 하는 방향으로 이루어짐

### ■ 변수 추출(variable/feature extraction)

- ✓ 원래의 변수 집단을 보다 효율적인 적은 수의 새로운 변수 집단으로 변환
- ✓ 데이터마이닝 모델에 독립적인 성능 지표가 추출된 변수의 효과를 측정하는 데 사용됨



# 차원축소: Dimensionality Reduction

## ❖ 차원 감소 기법 (cont')

- 변수 선택과 변수 추출 비교

$X_1$	$X_2$	$X_3$	...	$X_n$
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...

변수 선택

$X_1$	$X_5$	$X_8$
...	...	...
...	...	...
...	...	...
...	...	...
...	...	...

변수 추출

$Z_1$	$Z_2$	$Z_3$
...	...	...
...	...	...
...	...	...
...	...	...
...	...	...

$$Z_1 = X_1 + 0.2 * X_2$$

$$Z_2 = X_3 - 2 * X_5$$

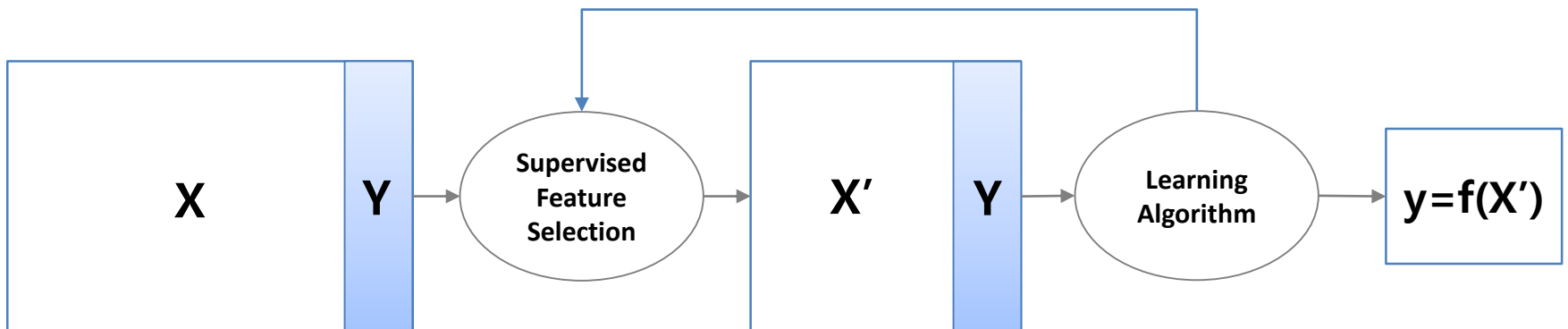
$$Z_3 = X_4 + X_6 - X_9$$



# 교사적 차원축소 기법

## ❖ 교사적 차원축소 기법 (Supervised feature selection)

- d-차원의 데이터에 대하여 사용하는 모델의 성능이 최대가 되도록 하는 d'차원( $d' \ll d$ )의 변수를 선택



- 변수 선택을 하기 전, 모델 구축에 사용할 알고리즘을 먼저 선택
- 동일한 데이터라도 모델 구축에 사용되는 알고리즘에 따라 다양한 선택 결과가 나타날 수 있음

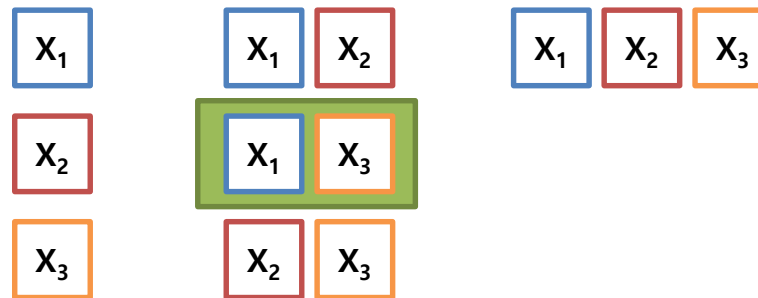
# Forward/Backward/Stepwise Selection

## ❖ 전역 탐색 (Exhaustive search)

- 가능한 모든 경우의 조합에 대해 모델을 구축한 뒤 최적의 변수 조합을 찾는 방식

✓ 예: 3개의 변수가 존재하는 경우  $x_1$   $x_2$   $x_3$

✓ 총 여섯 가지의 가능한 변수 조합 존재



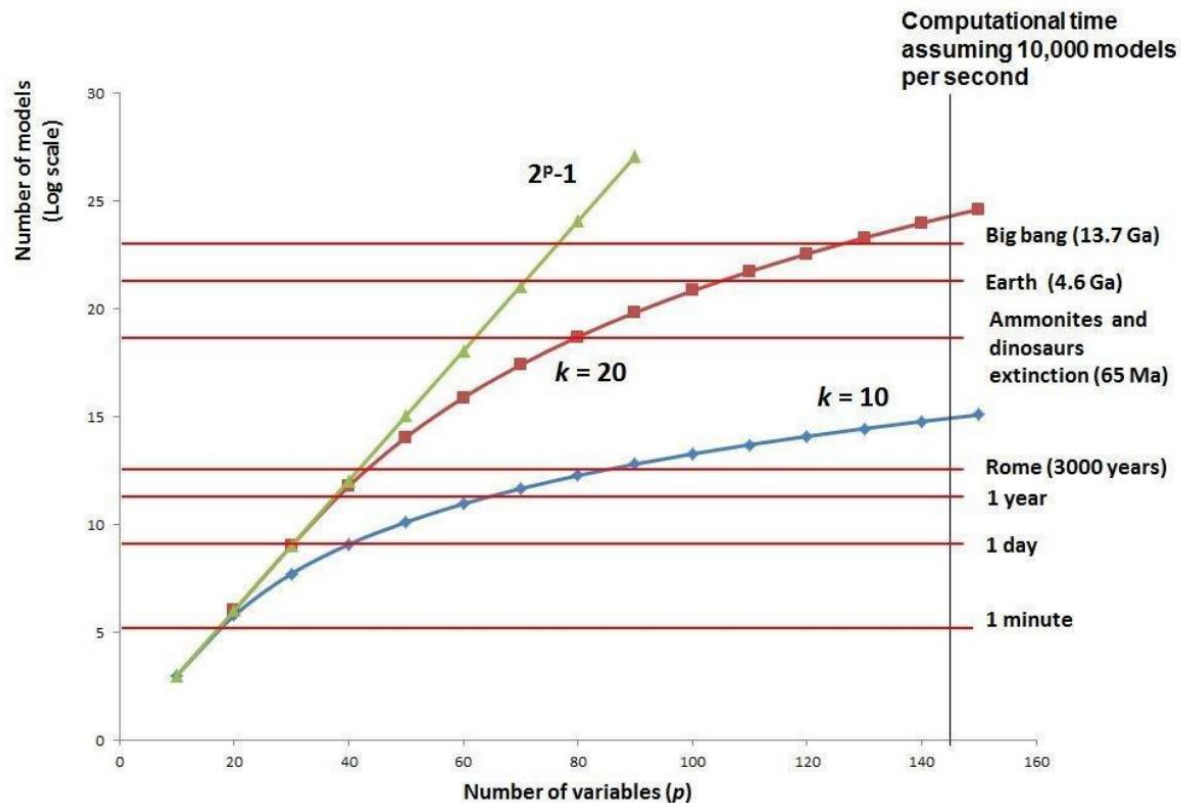
- 변수 선택을 위한 모델 평가 기준

✓ Akaike Information Criteria (AIC), Bayesian Information Criteria (BIC), 수정 R-제곱합, Mallows's  $C_p$  등

# Forward/Backward/Stepwise Selection

## ❖ 전역 탐색 (Exhaustive search)

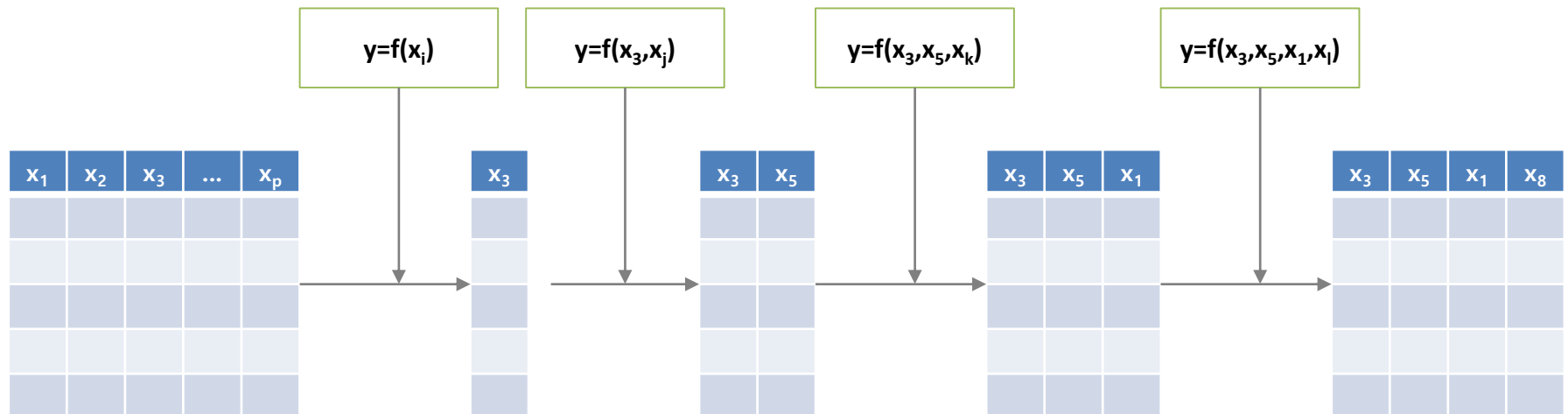
- 1초에 10,000개의 모델을 평가할 수 있는 컴퓨터를 활용할 경우 변수 선택에 소요되는 시간



# 전진 선택법 (Forward Selection)

## ❖ 전진 선택법

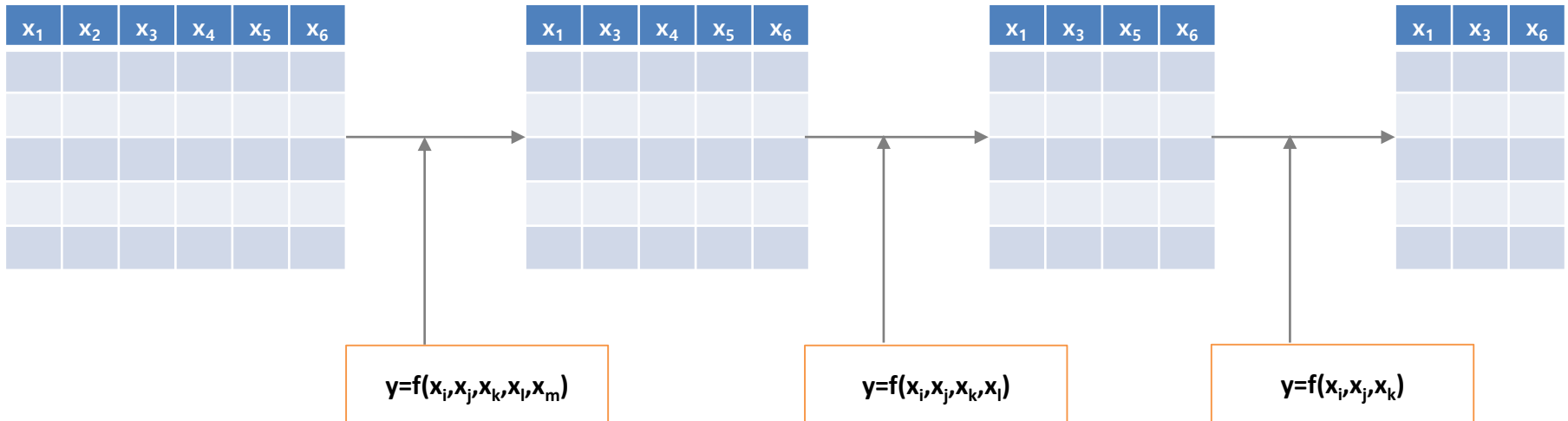
- 설명변수가 하나도 없는 모델에서부터 시작하여 가장 유의미한 변수를 하나씩 추가해 나가는 방법 (회귀분석 모델의 F-통계량 사용)
- 한번 선택된 변수는 제거되지 않음 (변수의 숫자는 단조 증가)
- 전진 선택법 예시



# 후진 소거법 (Backward Elimination)

## ❖ 후진 소거법

- 모든 변수를 사용하여 구축한 모델에서 유의미하지 않은 변수를 하나씩 제거해 나가는 방법
- 한번 제거된 변수는 다시 선택될 가능성이 없음 (변수의 숫자는 단조 증가)
- 후진 소거법 예시

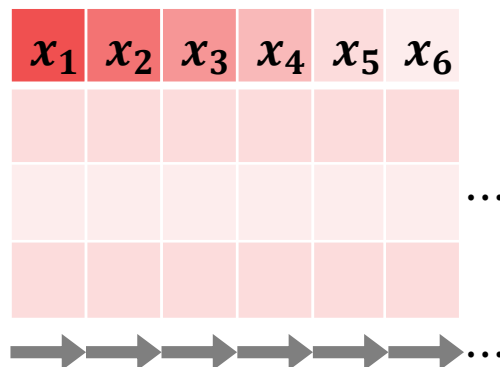


# 단계적 선택법 (Stepwise Selection)

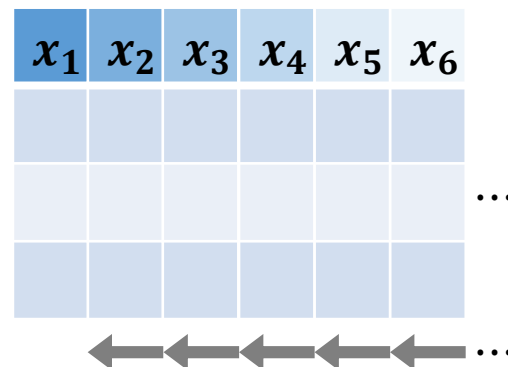
## ❖ 단계적 선택법

- 설명변수가 하나도 없는 모델에서부터 시작하여 전진선택법과 후진소거법을 번갈아가며 수행
- 전진선택법 및 후진소거법에 비해 시간을 오래 걸리나 보다 우수한 예측 성능을 나타내는 변수 집합을 찾아낼 가능성이 높음
- 한번 선택되거나 제거된 변수라도 다시 선택/제거될 가능성이 있음
- 변수의 수는 초기에는 일반적으로 증가하나 중반 이후에는 증가와 감소를 반복

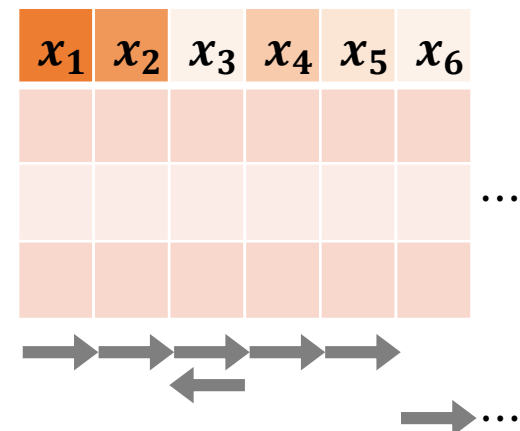
전진 선택법



후방 소거법

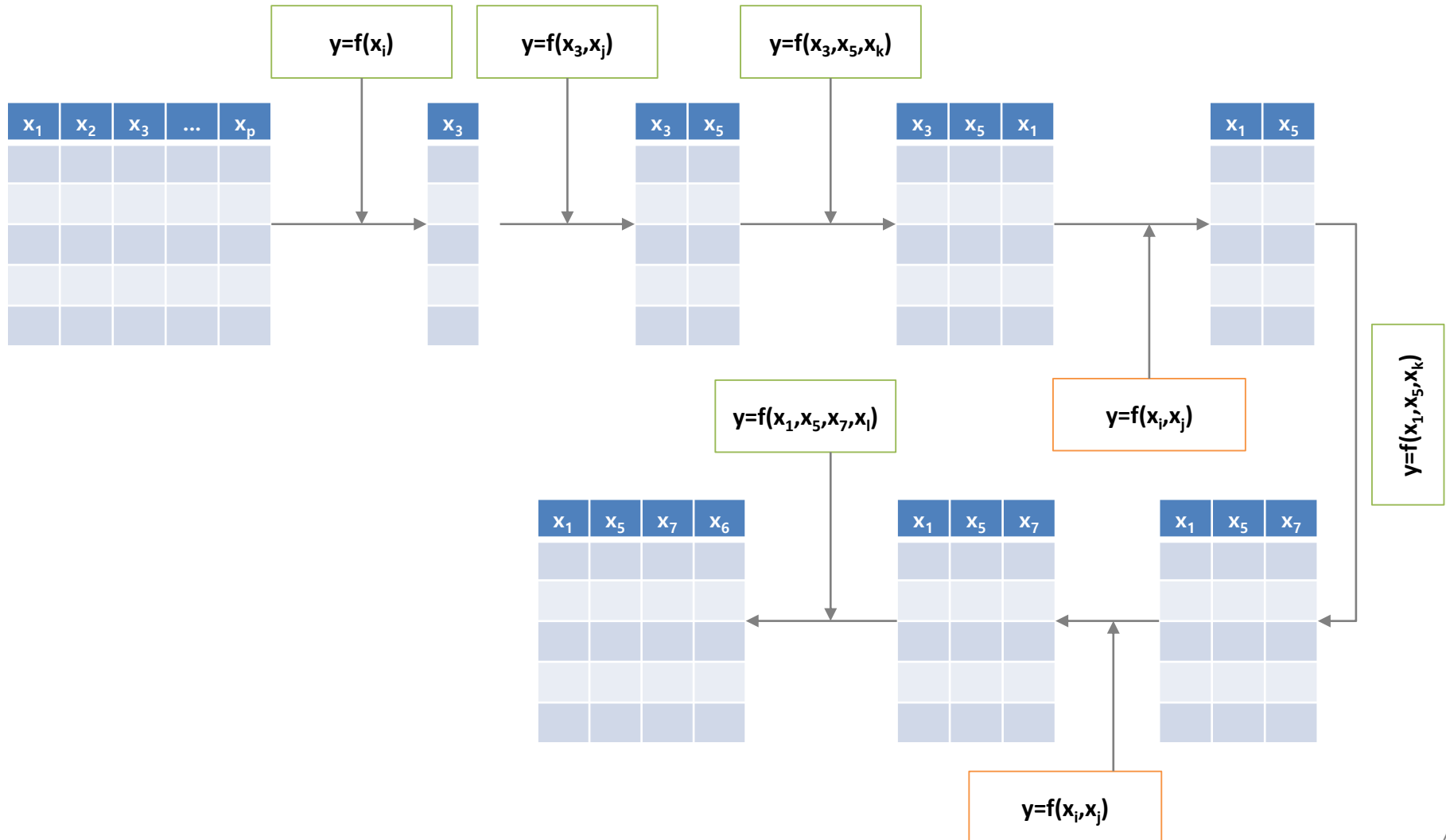


단계적 선택법



# 단계적 선택법 (Stepwise Selection)

## ❖ 단계적 선택법 예시



# 변수선택 평가지표

## ❖ 변수선택 평가지표 1 & 2

### ■ Akaike Information Criteria (AIC)

✓ 잔차제곱합에 변수의 수를 penalty term으로 추가

$$AIC = n \cdot \ln\left(\frac{SSE}{n}\right) + 2k$$

### ■ Bayesian Information Criteria (BIC)

✓ 잔차제곱합, 사용 변수의 수, 모든변수를 사용한 모델에서 추정된 잔차의 표준편차를 고려

$$BIC = n \cdot \ln\left(\frac{SSE}{n}\right) + \frac{2(k+2)n\sigma^2}{SSE} - \frac{2n^2\sigma^4}{SSE^2}$$



# 변수선택 평가지표

## ❖ 변수선택 평가지표 3

### ■ 수정 R-제곱합 (Adjusted $R^2$ )

✓ 단순 R-제곱합은 변수가 많아질수록 증가하므로 변수선택에 사용하기 좋은 평가지표가 아님

$$\text{Model 1 : } y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon$$

$$\text{Model 2 : } y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \dots + \beta_{k+m} x_{k+m} + \epsilon$$

$$R^2(M2) \geq R^2(M1)$$

✓ 변수의 수(k)를 고려한 수정 R-제곱합을 사용

$$\bar{R}^2 = 1 - \left( \frac{n-1}{n-k-1} \right) (1 - R^2) = 1 - \left( \frac{n-1}{n-k-1} \right) \frac{SSE}{SSTot}$$

# 변수선택 평가지표

## ❖ 변수선택 평가지표 4

### ■ Mallow's $C_p$

- ✓ 모델에 의해 설명되지 못하는 오차는 편기(Bias)와 분산(Variance)로 분해할 수 있음

$$\begin{aligned} \hat{y}_i - \mu_i &= (E[\hat{y}_i] - \mu_i) + (\hat{y}_i - E[\hat{y}_i]) & E[(\hat{y}_i - \mu_i)^2] &= (E[\hat{y}_i] - \mu_i)^2 + \text{Var}(\hat{y}_i) \\ &= \text{Bias} + \text{Random error} & &= \text{Bias}^2 + \text{Variance} \end{aligned}$$

- ✓ 다음과 같은 유도 과정을 통해 Mallow's  $C_p$ 를 최소화하는 변수 집합이 모델의 예측 성능을 최대화하는 것을 알 수 있음

$$\begin{aligned} \Gamma_p &= \frac{1}{\sigma^2} \left\{ \sum_{i=1}^n (E[\hat{y}_i] - \mu_i)^2 + \sum_{i=1}^n \text{Var}(\hat{y}_i) \right\} & \Gamma_p &= \frac{1}{\sigma^2} \{ E[SSE(p)] - (n-p)\sigma^2 + p\sigma^2 \} \\ &= \frac{SSB(p)}{\sigma^2} + \frac{1}{\sigma^2} \sum_{i=1}^n \text{Var}(\hat{y}_i) & &= \frac{E[SSE(p)]}{\sigma^2} - n + 2p \end{aligned}$$

p개의 선택된 변수를 이용한 모델의 SSE

$$C_p = \frac{SSE(p)}{MSE(K+1)} - n + 2p$$

모든 변수를 이용한 모델의 MSE

# 선형 회귀분석에서의 변수선택

## ❖ 전역 탐색 결과

#Coeffs	RSS	Cp	R <sup>2</sup>	Adj. R <sup>2</sup>	Prob	Model (Constant present in all models)											
						1	2	3	4	5	6	7	8	9	10	11	12
2	1,996,467,712	477.712	0.747	0.746	0.000	Constant	Age	*	*	*	*	*	*	*	*	*	*
3	1,672,546,432	305.506	0.788	0.787	0.000	Constant	Age	HP	*	*	*	*	*	*	*	*	*
4	1,438,242,432	181.495	0.818	0.817	0.000	Constant	Age	HP	Weight	*	*	*	*	*	*	*	*
5	1,258,062,976	86.594	0.840	0.839	0.000	Constant	Age	Mileage	HP	Weight	*	*	*	*	*	*	*
6	1,181,816,320	47.588	0.850	0.849	0.000	Constant	Age	Mileage	Petrol	Quarterly_Tax	Weight	*	*	*	*	*	*
7	1,095,153,024	2.980	0.861	0.860	0.962	Constant	Age	Mileage	Petrol	HP	Quarterly_Tax	Weight	*	*	*	*	*
8	1,093,753,344	4.227	0.861	0.860	0.994	Constant	Age	Mileage	Petrol	HP	Automatic	Quarterly_Tax	Weight	*	*	*	*
9	1,093,557,120	6.122	0.861	0.859	0.989	Constant	Age	Mileage	Petrol	HP	Metalic_Color	Automatic	Quarterly_Tax	Weight	*	*	*
10	1,093,422,592	8.049	0.861	0.859	0.976	Constant	Age	Mileage	Diesel	Petrol	HP	Metalic_Color	Automatic	Quarterly_Tax	Weight	*	*
11	1,093,335,424	10.002	0.861	0.859	0.961	Constant	Age	Mileage	Diesel	Petrol	HP	Metalic_Color	Automatic	CC	Quarterly_Tax	Weight	*
12	1,093,331,072	12.000	0.861	0.859	1.000	Constant	Age	Mileage	Diesel	Petrol	HP	Metalic_Color	Automatic	CC	Doors	Quarterly_Tax	Weight

# 선형 회귀분석에서의 변수선택

## ❖ 후진소거법 결과

#Coeffs	RSS	Cp	R <sup>2</sup>	Adj. R <sup>2</sup>	Prob	Model (Constant present in all models)											
						1	2	3	4	5	6	7	8	9	10	11	12
2	1,996,467,712	477.712	0.747	0.746	0.000	Constant	Age	*	*	*	*	*	*	*	*	*	*
3	1,780,184,064	363.394	0.774	0.773	0.000	Constant	Age	Weight	*	*	*	*	*	*	*	*	*
4	1,482,806,272	205.462	0.812	0.811	0.000	Constant	Age	Petrol	Weight	*	*	*	*	*	*	*	*
5	1,310,214,400	114.641	0.834	0.833	0.000	Constant	Age	Petrol	Quarterly_Tax	Weight	*	*	*	*	*	*	*
6	1,181,816,320	47.588	0.850	0.849	0.000	Constant	Age	Mileage	Petrol	Quarterly_Tax	Weight	*	*	*	*	*	*
7	1,095,153,024	2.980	0.861	0.860	0.962	Constant	Age	Mileage	Petrol	HP	Quarterly_Tax	Weight	*	*	*	*	*
8	1,093,753,344	4.227	0.861	0.860	0.994	Constant	Age	Mileage	Petrol	HP	Automatic	Quarterly_Tax	Weight	*	*	*	*
9	1,093,557,120	6.122	0.861	0.859	0.989	Constant	Age	Mileage	Petrol	HP	Metalic_Color	Automatic	Quarterly_Tax	Weight	*	*	*
10	1,093,422,592	8.049	0.861	0.859	0.976	Constant	Age	Mileage	Diesel	Petrol	HP	Metalic_Color	Automatic	Quarterly_Tax	Weight	*	*
11	1,093,335,424	10.002	0.861	0.859	0.961	Constant	Age	Mileage	Diesel	Petrol	HP	Metalic_Color	Automatic	CC	Quarterly_Tax	Weight	*
12	1,093,331,072	12.000	0.861	0.859	1.000	Constant	Age	Mileage	Diesel	Petrol	HP	Metalic_Color	Automatic	CC	Doors	Quarterly_Tax	Weight

# 선형 회귀분석에서의 변수선택

## ❖ 선택된 변수를 이용한 회귀분석 모델

### The Regression Model

Input variables	Coefficient	Std. Error	p-value	SS
Constant term	-3874.492188	1415.003052	0.00640071	97276411904
Age_08_04	-123.4366303	3.33806777	0	8033339392
KM	-0.01749926	0.00173714	0	251574528
Fuel_Type_Petrol	2409.154297	319.5795288	0	5049567
HP	19.70204735	4.22180223	0.00000394	291336576
Quarterly_Tax	16.88731384	2.08484554	0	192390864
Weight	15.91809368	1.26474357	0	281026176

### Training Data scoring - Summary Report

### Training Data scoring - Summary Report

Model Fit

Total sum of squared errors	RMS Error	Average Error
1516825972	1326.521353	-0.000143957

Total sum of squared errors	RMS Error	Average Error
1514553377	1325.527246	-0.000426154

### Validation Data scoring - Summary Report

### Validation Data scoring - Summary Report

Predictive performance

(compare to 12-predictor model!)

Total sum of squared errors	RMS Error	Average Error
1021510219	1334.029433	118.4483556

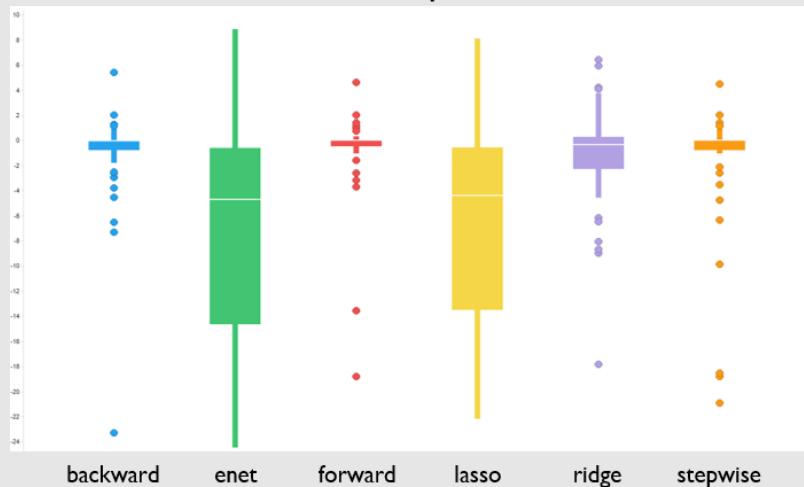
Total sum of squared errors	RMS Error	Average Error
1021587500	1334.079894	116.3728779

# 변수 선택 방식 실험

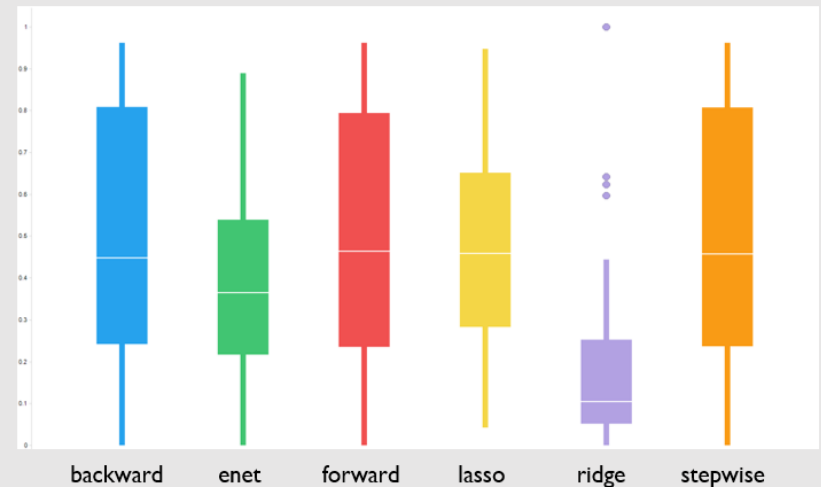
## ❖ 50개의 회귀분석 데이터

- 전진 선택, 후방 소거, 단계적 선택, Ridge, Lasso, Elastic Net

error rate improvement



변수 감소 비율



변수선택 방법	예측 정확도	변수 감소율	계산 효율성
Forward	4	4	1
Backward	3	3	2
Stepwise	2	2	6
Ridge	1	6	5
Lasso	6	1	3
Elastic Net	5	5	4

# 목차

I

다중선형회귀분석

II

회귀분석 성능 평가

III

변수 선택

IV

R 실습

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 선형회귀분석 데이터: Toyota Corolla 중고차가격 예측



Variable	Description	Variable	Description
		Guarantee_Period	Guarantee period in months
		ABS	Anti-Lock Brake System (Yes=1, No=0)
Price	Offer Price in EUROS	Airbag_1	Driver_Airbag (Yes=1, No=0)
Age_08_04	Age in months as in August 2004	Airbag_2	Passenger Airbag (Yes=1, No=0)
Mfg_Month	Manufacturing month (1-12)	Airco	Airconditioning (Yes=1, No=0)
Mfg_Year	Manufacturing Year	Automatic_airco	Automatic Airconditioning (Yes=1, No=0)
KM	Accumulated Kilometers on odometer	Boardcomputer	Boardcomputer (Yes=1, No=0)
Fuel_Type	Fuel Type (Petrol, Diesel, CNG)	CD_Player	CD Player (Yes=1, No=0)
HP	Horse Power	Central_Lock	Central Lock (Yes=1, No=0)
Met_Color	Metallic Color? (Yes=1, No=0)	Powered_Windows	Powered Windows (Yes=1, No=0)
Automatic	Automatic (Yes=1, No=0)	Power_Steering	Power Steering (Yes=1, No=0)
CC	Cylinder Volume in cubic centimeters	Radio	Radio (Yes=1, No=0)
Doors	Number of doors	Mistlamps	Mistlamps (Yes=1, No=0)
Cylinders	Number of cylinders	Sport_Model	Sport Model (Yes=1, No=0)
Gears	Number of gear positions	Backseat_Divider	Backseat Divider (Yes=1, No=0)
Quarterly_Tax	Quarterly road tax in EUROS	Metallic_Rim	Metallic Rim (Yes=1, No=0)
Weight	Weight in Kilograms	Radio_cassette	Radio Cassette (Yes=1, No=0)
Mfr_Guarantee	Within Manufacturer's Guarantee period (Yes=1, No=0)	Parking_Assistant	Parking assistance system (Yes=1, No=0)
BOVAG_Guarantee	BOVAG (Dutch dealer network) Guarantee (Yes=1, No=0)	Tow_Bar	Tow Bar (Yes=1, No=0)

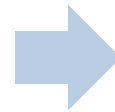


# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 데이터 불러오기 & 전처리

- 범주형 변수(선형회귀분석 사용 불가능)를 이진형 변수로 변환

Price	Age_08_04	Mfg_Month	Mfg_Year	KM	Fuel_Type	HP	Met_Color	Automatic	cc
13500	23	10	2002	46986	Diesel	90	1	0	2000
13750	23	10	2002	72937	Diesel	90	1	0	2000
13950	24	9	2002	41711	Diesel	90	1	0	2000
14950	26	7	2002	48000	Diesel	90	0	0	2000
13750	30	3	2002	38500	Diesel	90	0	0	2000
12950	32	1	2002	61000	Diesel	90	0	0	2000
16900	27	6	2002	94612	Diesel	90	1	0	2000
18600	30	3	2002	75889	Diesel	90	1	0	2000
21500	27	6	2002	19700	Petrol	192	0	0	1800
12950	23	10	2002	71138	Diesel	69	0	0	1900
20950	25	8	2002	31461	Petrol	192	0	0	1800
19950	22	11	2002	43610	Petrol	192	0	0	1800
19600	25	8	2002	32189	Petrol	192	0	0	1800
21500	31	2	2002	23000	Petrol	192	1	0	1800
22500	32	1	2002	34131	Petrol	192	1	0	1800



KM	HP	Met_Color
46986	90	1
72937	90	1
41711	90	1
48000	90	0
38500	90	0
61000	90	0
94612	90	1
75889	90	1
19700	192	0
71138	69	0
31461	192	0
43610	192	0
32189	192	0
23000	192	1

...

Petrol	Diesel	CNG
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
1	0	0
0	1	0
1	0	0
1	0	0
1	0	0
1	0	0
1	0	0

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 데이터 불러오기

```
# Read data from a file
corolla <- read.csv("ToyotaCorolla.csv")

# Indices for the activated input variables
nCar <- dim(corolla)[1]
nVar <- dim(corolla)[2]

id_idx <- c(1,2)
category_idx <- 8
```

### ■ read.csv( ): csv 파일을 읽어들이는 함수

- ✓ Full path가 지정되어 있지 않으면 rstudio에서 지정한 working directory에서 해당 이름을 가진 파일을 찾아서 읽어들이м
- ✓ dim( ): data.frame의 행과 열을 반환하는 함수
- ✓ id\_idx: id와 관련하여 분석이 필요없는 열 번호
- ✓ category\_idx: 명목형 변수로서 1-of-C coding 변환이 필요한 열 번호

# R 실습: 다중선택회귀분석 및 변수선택

## ❖ 명목형 변수의 1-of-C coding 처리

```
# 범주형 변수를 이진형 변수로 변환
dummy_p <- rep(0, nCar)
dummy_d <- rep(0, nCar)
dummy_c <- rep(0, nCar)

p_idx <- which(corolla$Fuel_Type == "Petrol")
d_idx <- which(corolla$Fuel_Type == "Diesel")
c_idx <- which(corolla$Fuel_Type == "CNG")

dummy_p[p_idx] <- 1
dummy_d[d_idx] <- 1
dummy_c[c_idx] <- 1
```

- `dummy_p(c/d)`: 총 `nCar`개의 길이를 갖는 벡터의 모든 값을 0으로 초기화
- `p_idx`: `corolla` 데이터의 `Fule_Type`이 “Petrol”인 행의 번호를 저장 (`d_idx`, `c_idx`도 유사한 작업 수행)
- `dummy_p[p_idx] <- 1`: `dummy_p` 벡터에 `p_idx`에 해당하는 원소들의 값을 1로 대체 (`d_idx`, `c_idx`도 유사한 작업 수행)

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 명목형 변수의 1-of-C coding 처리 및 데이터 결합

```
Fuel <- data.frame(dummy_p, dummy_d, dummy_c)
names(Fuel) <- c("Petrol", "Diesel", "CNG")

# Prepare the data for MLR
mlr_data <- cbind(corolla[, -c(id_idx, category_idx)], Fuel)
```

- dummy\_p, dummy\_d, dummy\_c를 결합하여 Fuel이라는 data.frame 생성
- Fuel 데이터프레임의 각 열에 “Petrol”, “Diesel”, “CNG” 이름 부여
- 원래 corolla 데이터에서 id\_idx와 category\_idx에 해당하는 열을 제외하고 새로 생성한 Fuel 데이터를 결합(cbind( )함수 사용)

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 학습용과 검증용 데이터 분할

```
# Split the data into the training/validation sets
set.seed(12345)
trn_idx <- sample(1:nCar, round(0.7*nCar))
trn_data <- mlr_data[trn_idx,]
val_data <- mlr_data[-trn_idx,]
```

- `set.seed(12345)`: 난수 생성을 위한 설정을 모두 동일하게 → 모든 사람의 학습 데이터가 같고 모든 사람의 검증용 데이터가 같아짐
- `trn_idx`: `sample( )` 함수(1부터 `nCar`까지의 자연수 중에서 `nCar`의 70%에 해당하는 수만큼을 비복원 랜덤 추출)를 사용하여 학습 데이터 행번호 저장
- `trn_data`: 학습 데이터
- `val_data`: 검증 데이터(학습 데이터 행번호에 해당하지 않는 행번호들로 구성)

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 모든 변수를 이용한 모델 학습

```
# Train the MLR
full_model <- lm(Price ~ ., data = trn_data)
full_model
summary(full_model)
plot(full_model)
```

- `lm()`: linear model을 생성해주는 함수
  - ✓ `Price ~ .`: Formula 양식 (~ 왼쪽은 종속변수, 오른쪽은 설명변수를 의미, 두 개 이상의 설명변수는 +기호로 결합, 종속변수 이외의 모든 변수를 설명변수로 사용할 때는 `period(.)` 사용)
  - ✓ `data = trn_data`: 앞서 생성한 `trn_data`를 이용하여 회귀계수 추정
- `summary()`: 회귀계수 추정이 완료된 모형에 대한 여러 정보를 콘솔창에 출력
- `plot()`: 선형회귀분석 관련 그래프 도시

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 모든 변수를 이용한 모델 학습

### ■ 결과 해석

- ✓ Estimate: 추정된 회귀 계수
- ✓ Std.Error: 추정된 회귀 계수의 표준편차
- ✓ t value: 해당 변수의 유의미함에 대한 통계량
- ✓ Pr(>|t|): 변수의 유의 확률 (0에 가까울수록 유의미한 변수)
- ✓ Adjusted R-squared: 수정 R제곱합 (1에 가까울수록 회귀모형이 종속변수를 잘 설명함을 의미)
- ✓ NA: 해당 변수는 다중공선성 문제로 모형에서 제외됨

```
> summary(full_model)
```

```
Call:
lm(formula = Price ~ ., data = trn_data)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-8569.6  -637.3   -42.9   650.5  5720.8
```

```
Coefficients: (3 not defined because of singularities)
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -752.65641  1775.96778  -0.424  0.671805
Age_08_04    -118.43999    4.29358  -27.585 < 2e-16 ***
Mfg_Month     -95.78658    11.11831   -8.615 < 2e-16 ***
Mfg_Year              NA              NA      NA      NA
KM             -0.01727    0.00136  -12.702 < 2e-16 ***
HP              20.46048    3.59449    5.692  1.66e-08 ***
Met_Color     -64.93066    81.74804   -0.794  0.427228
Automatic      338.02084   156.47529    2.160  0.031000 *
cc             -0.10770    0.07958   -1.353  0.176246
Doors          10.46213    43.60079    0.240  0.810418
Cylinders              NA              NA      NA      NA
Gears         183.36459   196.28703    0.934  0.350451
```

•  
•  
•

```
Tow_Bar      -217.67837    85.11397  -2.557  0.010694 *
Petrol        2280.57458   387.15749    5.891  5.30e-09 ***
Diesel       1004.02078   377.75296    2.658  0.007993 **
CNG              NA              NA      NA      NA
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

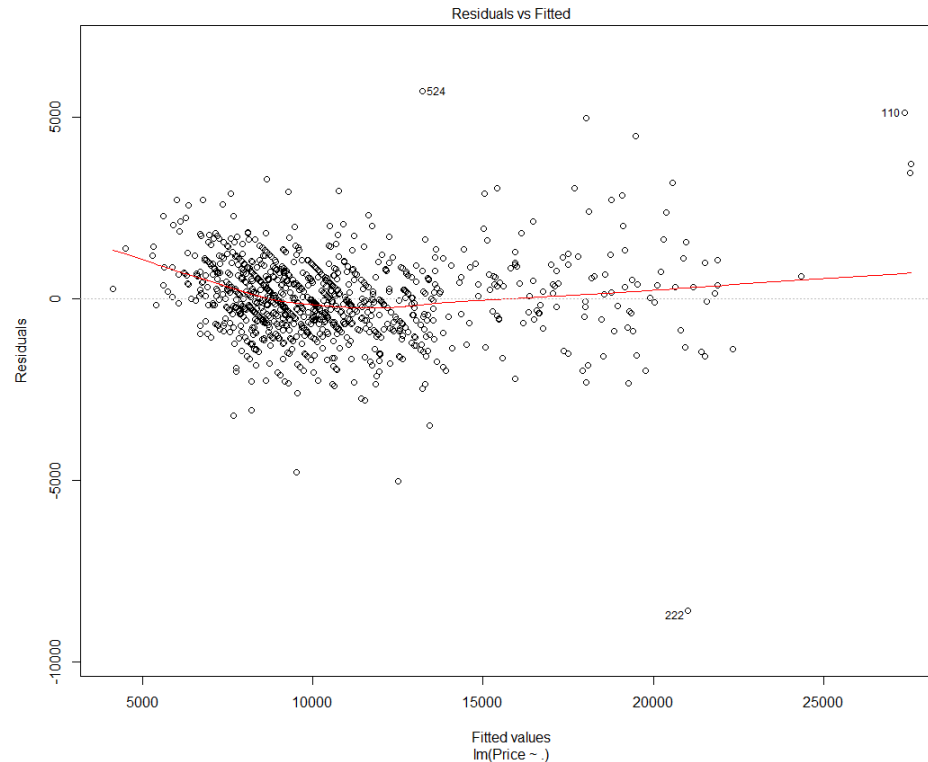
```
Residual standard error: 1128 on 971 degrees of freedom
Multiple R-squared:  0.9046,    Adjusted R-squared:  0.9014
F-statistic: 279.1 on 33 and 971 DF,  p-value: < 2.2e-16
```

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 모든 변수를 이용한 모델 학습

### ■ 그림 1: 아래 가정을 검증하는데 사용

- ✓ 종속변수 Y에 대한 오차항(residual)은 설명변수 값의 범위에 관계없이 일정함(homoskedasticity)



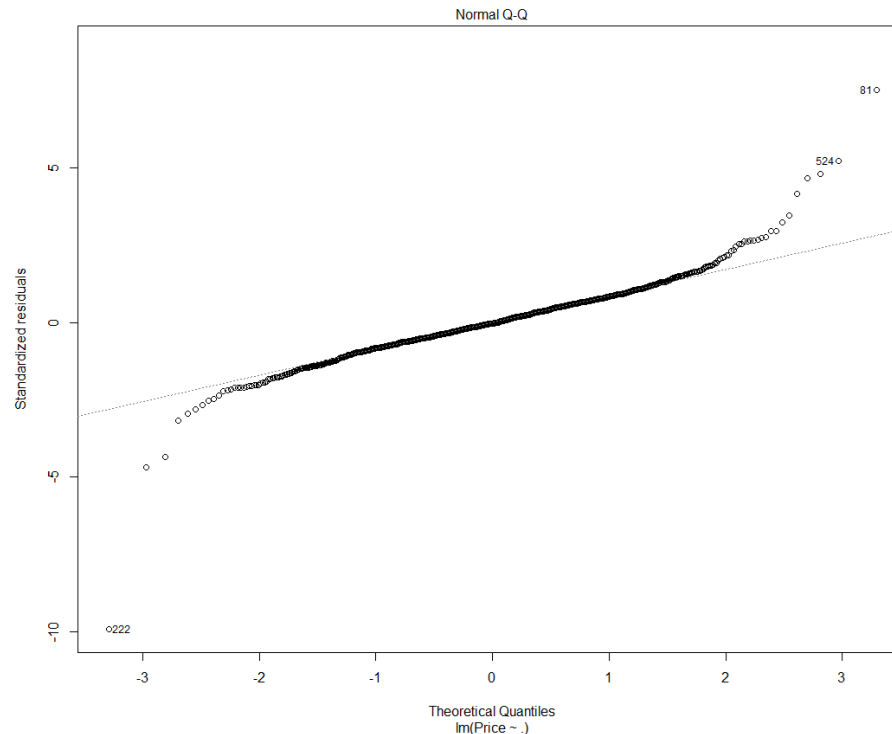


# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 모든 변수를 이용한 모델 학습

### ■ 그림 2: 아래 가정을 검증하는데 사용

- 오차항  $\varepsilon$  이 정규분포를 따름
- ✓ x축 기준 +2 초과 -2 미만 부분에서 점들이 직선 위로부터 벗어나면 어느 정도 가정 만족하는 것으로 판단

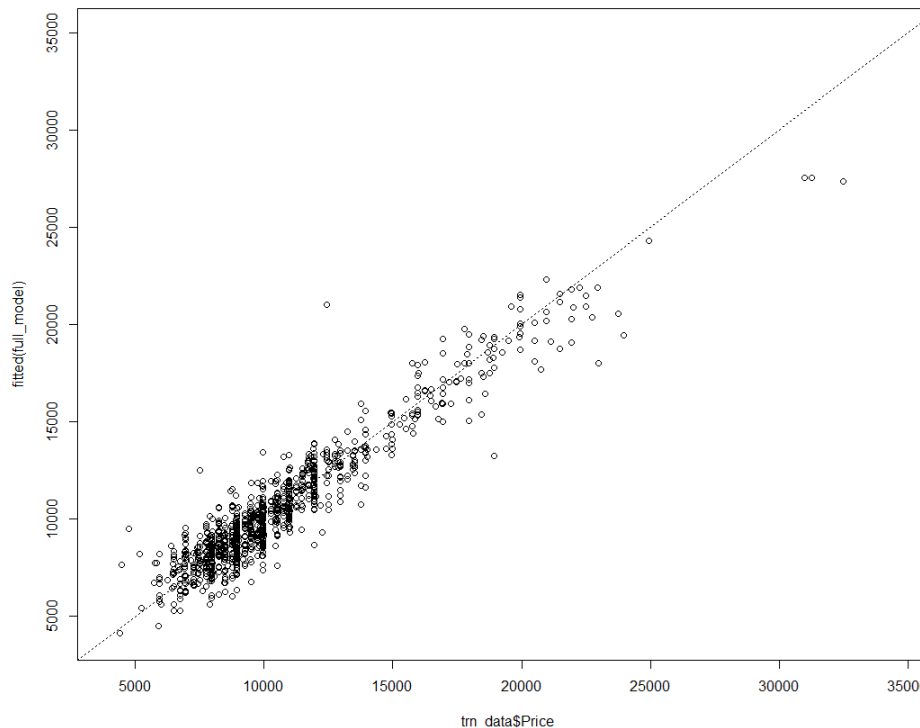


# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 모든 변수를 이용한 모델 학습

```
# Plot the result  
plot(trn_data$Price, fitted(full_model), xlim = c(4000,35000), ylim =  
c(4000,35000))  
abline(0,1,lty=3)
```

- 실제 값(x축)과 추정된 값(y축) 사이의 관계 도시



# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 전진선택법을 이용한 변수 선택

```
# 변수선택 1: 전진선택법
# Upperbound formula 만들기
tmp_x <- paste(colnames(trn_data)[-1], collapse=" + ")
tmp_xy <- paste("Price ~ ", tmp_x, collapse = "")
as.formula(tmp_xy)
```

- 선형회귀분석 수행 시 Domain knowledge를 가지고 반드시 포함되어야 하는 변수 또는 반드시 제거되어야 하는 변수를 미리 설정할 수 있음

- ✓ 변수 집합의 상한(upperbound)을 정의할 때, 반드시 제거되어야 하는 변수를 포함시키지 않도록 설정

- ✓ 본 실습에서는 모든 변수가 포함되도록 설정

```
> as.formula(tmp_xy)
```

```
Price ~ Age_08_04 + Mfg_Month + Mfg_Year + KM + HP + Met_Color +
Automatic + cc + Doors + Cylinders + Gears + Quarterly_Tax +
Weight + Mfr_Guarantee + BOVAG_Guarantee + Guarantee_Period +
ABS + Airbag_1 + Airbag_2 + Airco + Automatic_airco + Boardcomputer +
CD_Player + Central_Lock + Powered_Windows + Power_Steering +
Radio + Mistlamps + Sport_Model + Backseat_Divider + Metallic_Rim +
Radio cassette + Tow Bar + Petrol + Diesel + CNG
```

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 전진선택법을 이용한 변수 선택

```
forward_model <- step(lm(Price ~ 1, data = trn_data),  
                      scope = list(upper = as.formula(tmp_xy), lower = Price ~ 1),  
                      direction="forward", trace=1)  
summary(forward_model)
```

- `step()`: 전진선택/후방소거/단계적선택법 수행 함수
  - ✓ 첫 번째 인수: 시작 모형 (여기서는 아무 변수도 사용하지 않은 모형)
  - ✓ 두 번째 인수: 모형의 탐색 범위 (여기서는 아무 변수를 사용하지 않는 것부터 모든 변수를 사용하는 것까지)
  - ✓ 세 번째 인수(`direction = “ ”`): 변수선택 방향을 설정 (“forward”가 전진 선택)
  - ✓ 네 번째 인수: 변수 선택 과정을 콘솔창에 출력

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 전진선택법 결과

	변수의 수	수정 R-제곱합
모든변수	36	0.9014
전진선택법	22	0.9018
후방소거법		
단계적선택법		

- 변수 12개 감소
- 수정 R-제곱합 0.0004 증가

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-2.870e+06	9.915e+04	-28.949	< 2e-16 ***
Mfg_Year	1.432e+03	4.975e+01	28.776	< 2e-16 ***
Automatic_airco	2.216e+03	1.975e+02	11.222	< 2e-16 ***
KM	-1.717e-02	1.341e-03	-12.803	< 2e-16 ***
Weight	1.205e+01	1.296e+00	9.299	< 2e-16 ***
HP	2.060e+01	3.424e+00	6.018	2.49e-09 ***
Guarantee_Period	7.961e+01	1.500e+01	5.309	1.36e-07 ***
BOVAG_Guarantee	4.878e+02	1.351e+02	3.611	0.00032 ***
Powered_Windows	3.223e+02	8.139e+01	3.960	8.05e-05 ***
Quarterly_Tax	1.517e+01	1.849e+00	8.208	7.03e-16 ***
Petrol	2.314e+03	3.759e+02	6.156	1.09e-09 ***
Tow_Bar	-2.181e+02	8.284e+01	-2.633	0.00859 **
Metallic_Rim	2.415e+02	9.668e+01	2.498	0.01267 *
CD_Player	2.482e+02	1.067e+02	2.327	0.02015 *
Backseat_Divider	-3.576e+02	1.238e+02	-2.888	0.00397 **
Sport_Model	2.364e+02	9.092e+01	2.600	0.00947 **
Mfr_Guarantee	1.956e+02	7.801e+01	2.507	0.01233 *
Diesel	9.928e+02	3.677e+02	2.700	0.00705 **
Automatic	2.889e+02	1.530e+02	1.889	0.05924 .
Boardcomputer	-2.784e+02	1.272e+02	-2.189	0.02884 *
ABS	-2.111e+02	1.048e+02	-2.015	0.04423 *
Mfg_Month	2.344e+01	1.098e+01	2.135	0.03297 *
Radio_cassette	-2.059e+02	1.083e+02	-1.902	0.05751 .

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

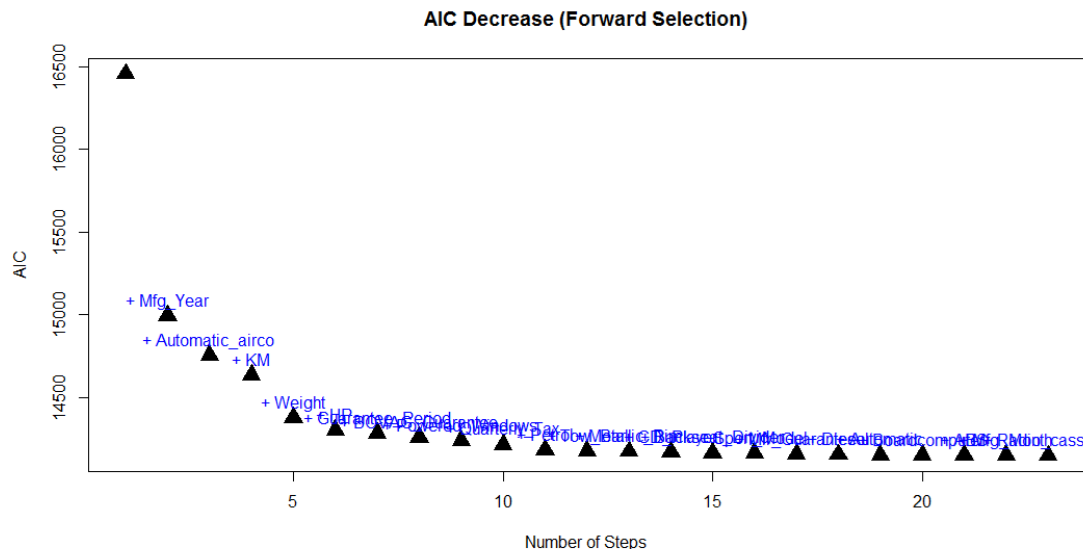
Residual standard error: 1126 on 982 degrees of freedom  
 Multiple R-squared: 0.904, Adjusted R-squared: 0.9018  
 F-statistic: 420.2 on 22 and 982 DF, p-value: < 2.2e-16

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 전진선택법을 이용한 변수 선택

```
# 각 단계에서 선택된 변수 표시
forward_model$anova$Step
forward_model$anova$AIC
# 선택된 변수에 따른 AIC 감소분 표시
plot(forward_model$anova$AIC, pch = 17, cex=2,
      main = "AIC Decrease (Forward Selection)",
      xlab = "Number of Steps", ylab = "AIC")
text(forward_model$anova$AIC, forward_model$anova$Step, cex=1, pos=3, col="blue")
```

### ■ 각 단계에서 선택된 변수와 AIC 감소분에 대한 정보



# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 후방선택법을 이용한 변수 선택

```
# 변수선택 2: 후방소거법
backward_model <- step(full_model, scope = list(upper = as.formula(tmp_xy),
  lower = Price ~ 1), direction="backward", trace=1)
summary(backward_model)
```

- `step()`: 전진선택/후방소거/단계적선택법 수행 함수
  - ✓ 첫 번째 인수: 시작 모형 (여기서는 모든 변수를 사용한 모형(`full_model`))
  - ✓ 두 번째 인수: 모형의 탐색 범위 (여기서는 아무 변수를 사용하지 않는 것부터 모든 변수를 사용하는 것까지)
  - ✓ 세 번째 인수(`direction = “ ”`): 변수선택 방향을 설정 (“backward”가 후방 소거)
  - ✓ 네 번째 인수: 변수 선택 과정을 콘솔창에 출력

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 후방소거법 결과

	변수의 수	수정 R-제곱합
모든변수	36	0.9014
전진선택법	22	0.9018
후방소거법	22	0.9018
단계적선택법		

- 변수 12개 감소
- 수정 R-제곱합 0.0004 증가
- 본 실습은 전진선택법과 후방소거법에 의해 선택된 변수는 일치하나 계수가 다름

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-7.051e+01	1.400e+03	-0.050	0.95983
Age_08_04	-1.193e+02	4.146e+00	-28.776	< 2e-16 ***
Mfg_Month	-9.587e+01	1.102e+01	-8.700	< 2e-16 ***
KM	-1.717e-02	1.341e-03	-12.803	< 2e-16 ***
HP	2.060e+01	3.424e+00	6.018	2.49e-09 ***
Automatic	2.889e+02	1.530e+02	1.889	0.05924 .
Quarterly_Tax	1.517e+01	1.849e+00	8.208	7.03e-16 ***
Weight	1.205e+01	1.296e+00	9.299	< 2e-16 ***
Mfr_Guarantee	1.956e+02	7.801e+01	2.507	0.01233 *
BOVAG_Guarantee	4.878e+02	1.351e+02	3.611	0.00032 ***
Guarantee_Period	7.961e+01	1.500e+01	5.309	1.36e-07 ***
ABS	-2.111e+02	1.048e+02	-2.015	0.04423 *
Automatic_airco	2.216e+03	1.975e+02	11.222	< 2e-16 ***
Boardcomputer	-2.784e+02	1.272e+02	-2.189	0.02884 *
CD_Player	2.482e+02	1.067e+02	2.327	0.02015 *
Powered_Windows	3.223e+02	8.139e+01	3.960	8.05e-05 ***
Sport_Model	2.364e+02	9.092e+01	2.600	0.00947 **
Backseat_Divider	-3.576e+02	1.238e+02	-2.888	0.00397 **
Metallic_Rim	2.415e+02	9.668e+01	2.498	0.01267 *
Radio_cassette	-2.059e+02	1.083e+02	-1.902	0.05751 .
Tow_Bar	-2.181e+02	8.284e+01	-2.633	0.00859 **
Petrol	2.314e+03	3.759e+02	6.156	1.09e-09 ***
Diesel	9.928e+02	3.677e+02	2.700	0.00705 **

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1126 on 982 degrees of freedom  
Multiple R-squared: 0.904, Adjusted R-squared: 0.9018  
F-statistic: 420.2 on 22 and 982 DF, p-value: < 2.2e-16

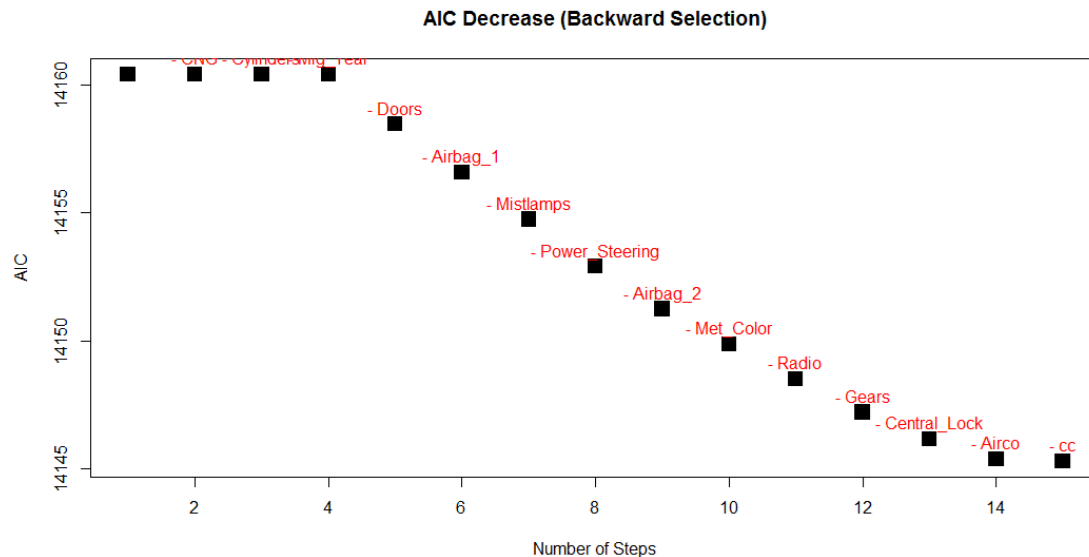


# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 후방소거법을 이용한 변수 선택

```
# 각 단계에서 제거된 변수 표시
backward_model$anova$Step
backward_model$anova$AIC
# 제거된 변수에 따른 AIC 감소분 표시
plot(backward_model$anova$AIC, pch = 15, cex=2,
      main = "AIC Decrease (Backward Selection)",
      xlab = "Number of Steps", ylab = "AIC")
text(backward_model$anova$AIC, backward_model$anova$Step, cex=1, pos=3, col="red")
```

### ■ 각 단계에서 제거된 변수와 AIC 감소분에 대한 정보



# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 단계적 선택법을 이용한 변수 선택

```
# 변수선택 3: 단계적 선택법
stepwise_model <- step(lm(Price ~ 1, data = trn_data),
  scope = list(upper = as.formula(tmp_xy), lower = Price ~ 1),
  direction="both", trace=1)
summary(stepwise_model)
```

- `step()`: 전진선택/후방소거/단계적선택법 수행 함수
  - ✓ 첫 번째 인수: 시작 모형 (여기서는 아무 변수도 사용하지 않은 모형)
  - ✓ 두 번째 인수: 모형의 탐색 범위 (여기서는 아무 변수를 사용하지 않는 것부터 모든 변수를 사용하는 것까지)
  - ✓ 세 번째 인수(`direction = “ ”`): 변수선택 방향을 설정 (“both”가 단계적 선택)
  - ✓ 네 번째 인수: 변수 선택 과정을 콘솔창에 출력

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 후방소거법 결과

	변수의 수	수정 R-제곱합
모든변수	36	0.9014
전진선택법	22	0.9018
후방소거법	22	0.9018
단계적선택법	22	0.9018

- 변수 12개 감소
- 수정 R-제곱합 0.0004 증가
- 전진선택법과 동일한 결과(선택된 변수 및 회귀계수가 모두 일치)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-2.870e+06	9.915e+04	-28.949	< 2e-16 ***
Mfg_Year	1.432e+03	4.975e+01	28.776	< 2e-16 ***
Automatic_airco	2.216e+03	1.975e+02	11.222	< 2e-16 ***
KM	-1.717e-02	1.341e-03	-12.803	< 2e-16 ***
Weight	1.205e+01	1.296e+00	9.299	< 2e-16 ***
HP	2.060e+01	3.424e+00	6.018	2.49e-09 ***
Guarantee_Period	7.961e+01	1.500e+01	5.309	1.36e-07 ***
BOVAG_Guarantee	4.878e+02	1.351e+02	3.611	0.00032 ***
Powered_Windows	3.223e+02	8.139e+01	3.960	8.05e-05 ***
Quarterly_Tax	1.517e+01	1.849e+00	8.208	7.03e-16 ***
Petrol	2.314e+03	3.759e+02	6.156	1.09e-09 ***
Tow_Bar	-2.181e+02	8.284e+01	-2.633	0.00859 **
Metallic_Rim	2.415e+02	9.668e+01	2.498	0.01267 *
CD_Player	2.482e+02	1.067e+02	2.327	0.02015 *
Backseat_Divider	-3.576e+02	1.238e+02	-2.888	0.00397 **
Sport_Model	2.364e+02	9.092e+01	2.600	0.00947 **
Mfr_Guarantee	1.956e+02	7.801e+01	2.507	0.01233 *
Diesel	9.928e+02	3.677e+02	2.700	0.00705 **
Automatic	2.889e+02	1.530e+02	1.889	0.05924 .
Boardcomputer	-2.784e+02	1.272e+02	-2.189	0.02884 *
ABS	-2.111e+02	1.048e+02	-2.015	0.04423 *
Mfg_Month	2.344e+01	1.098e+01	2.135	0.03297 *
Radio_cassette	-2.059e+02	1.083e+02	-1.902	0.05751 .

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1126 on 982 degrees of freedom  
Multiple R-squared: 0.904, Adjusted R-squared: 0.9018  
F-statistic: 420.2 on 22 and 982 DF, p-value: < 2.2e-16

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 네 가지 모형의 예측 정확도 비교

```
# 검증 데이터에 대한 각 변수선택 결과의 예측 정확도 비교
full_haty <- predict(full_model, newdata = val_data)
forward_haty <- predict(forward_model, newdata = val_data)
backward_haty <- predict(backward_model, newdata = val_data)
stepwise_haty <- predict(stepwise_model, newdata = val_data)
```

- `predict()`: 구축된 모형을 이용하여 새로운 데이터에 적용하여 예측 수행
  - ✓ 첫 번째 인수: 예측에 사용할 모형
  - ✓ 두 번째 인수: 예측 대상이 되는 새로운 데이터

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 네 가지 모형의 예측 정확도 비교

```
# 회귀분석 예측성능 평가지표
# 1: Mean squared error (MSE)
perf_mat <- matrix(0,4,4)
perf_mat[1,1] <- mean((val_data$Price-full_haty)^2)
perf_mat[1,2] <- mean((val_data$Price-forward_haty)^2)
perf_mat[1,3] <- mean((val_data$Price-backward_haty)^2)
perf_mat[1,4] <- mean((val_data$Price-stepwise_haty)^2)
# 2: Root mean squared error (RMSE)
perf_mat[2,1] <- sqrt(mean((val_data$Price-full_haty)^2))
perf_mat[2,2] <- sqrt(mean((val_data$Price-forward_haty)^2))
perf_mat[2,3] <- sqrt(mean((val_data$Price-backward_haty)^2))
perf_mat[2,4] <- sqrt(mean((val_data$Price-stepwise_haty)^2))
# 3: Mean absolute error (MAE)
perf_mat[3,1] <- mean(abs(val_data$Price-full_haty))
perf_mat[3,2] <- mean(abs(val_data$Price-forward_haty))
perf_mat[3,3] <- mean(abs(val_data$Price-backward_haty))
perf_mat[3,4] <- mean(abs(val_data$Price-stepwise_haty))
# 4: Mean absolute percentage error (MAPE)
perf_mat[4,1] <- mean(abs((val_data$Price-full_haty)/val_data$Price))*100
perf_mat[4,2] <- mean(abs((val_data$Price-forward_haty)/val_data$Price))*100
perf_mat[4,3] <- mean(abs((val_data$Price-backward_haty)/val_data$Price))*100
perf_mat[4,4] <- mean(abs((val_data$Price-stepwise_haty)/val_data$Price))*100
# 변수선택 기법 결과 비교
rownames(perf_mat) <- c("MSE", "RMSE", "MAE", "MAPE")
colnames(perf_mat) <- c("All", "Forward", "Backward", "Stepwise")
perf_mat
```

# R 실습: 다중선형회귀분석 및 변수선택

## ❖ 네 가지 모형의 예측 정확도 비교

```
> perf_mat
```

	All	Forward	Backward	Stepwise
MSE	1.177495e+06	1.179686e+06	1.179686e+06	1.179686e+06
RMSE	1.085124e+03	1.086133e+03	1.086133e+03	1.086133e+03
MAE	8.140598e+02	8.147627e+02	8.147627e+02	8.147627e+02
MAPE	8.007297e+00	8.002794e+00	8.002794e+00	8.002794e+00

- MSE, RMSE, MAE 관점에서는 모든 변수를 사용한 모형이 변수 선택을 수행한 모형에 비해서 우수함
- MAPE 관점에서는 변수 선택을 수행한 모형이 모든 변수를 사용한 모형에 비해서 우수함
- 변수 선택 기법 간의 성능 차이는 없음

# Q & A

